



یادگیری بازنمایی علی برای تعمیم خارج از توزیع

گزارش سمینار کارشناسی ارشد

در رشته مهندسی کامپیوتر - گرایش هوش مصنوعی و رباتیک

نام دانشجو:

حسین رضائی

استاد راهنما:

دکتر عادل رحمانی

آذر ماه 1401

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

چکیده

الگوریتم های یادگیری ماشین معمولاً روی فرض iid ساخته می شوند که داده های آموزش و آزمون مستقل و به صورت یکسان توزیع شده اند. این فرض به این معنا است که توزیع داده های آموزش و داده های آزمون یکسان باشند. در دنیا واقعی به دلیل تغییرات توزیعی، این فرض به سختی برآورده می شود که باعث کاهش شدید دقت این الگوریتم های کلاسیک یادگیری ماشین می شود. از طرف دیگر، الگوریتم های یادگیری ماشین غالباً از مدل های آماری برای مدل سازی وابستگی بین داده ها و برچسب ها استفاده می کنند که قصد یادگیری بازنمایی های مستقل از دامنه را دارد. با این وجود، مدل های آماری توصیف های سطحی واقعیت هستند، زیرا آنها فقط به مدل سازی وابستگی به جای مکانیسم علی ذاتی نیاز دارند. هنگامی که وابستگی با توزیع هدف تغییر می کند، مدل های آماری ممکن است در تعمیم ناکام باشند. علیت، با تمرکز بر بازنمایی دانش ساختاری در مورد فرآیند تولید داده که اجازه تغییرات (interventions) را می دهد، می تواند به درک و رفع برخی محدودیت های روش های یادگیری ماشین فعلی کمک کند. علیرغم موفقیت یادگیری آماری، این مدل ها توصیفی نسبتاً سطحی از واقعیت ارائه می دهند که تنها زمانی برقرار است که شرایط آزمایشی ثابت باشد. در عوض، حوزه یادگیری علی به دنبال مدل سازی اثر تغییرات توزیع با ترکیبی از یادگیری مبتنی بر داده و مفروضاتی است که قبلاً در توصیف آماری یک سیستم گنجانده نشده اند.

برای مشکل ناشی از تغییر توزیع، یعنی جایی که توزیع داده های آزمون متفاوت از داده های آموزش است، مسئله ی تعمیم خارج از توزیع مطرح می شود که در آن الگوریتم بتواند به خوبی عمل کند و تعمیم خوبی روی داده های دیده نشده یعنی داده های آزمون داشته باشد. در این سمینار ما به بررسی روش های مسئله ی تعمیم خارج از توزیع می پردازیم و همینطور یکی از روش های آن به نام علیت که اخیراً به صورت گستره مورد توجه قرار گرفته است را با جزئیات دقیق تر بررسی می کنیم.

واژه های کلیدی: تعمیم خارج از توزیع، استنتاج علی، یادگیری ثابت، یادگیری بازنمایی، تعمیم دامنه

فهرست مطالب

عنوان	صفحه
فصل 1: مقدمه	1
1-1- شرح مسأله.....	2
1-2- معرفی حوزه سمینار.....	3
1-3- ساختار گزارش.....	4
فصل 2: تعاریف و مفاهیم مبنایی	5
2-1- مقدمه.....	6
2-2- تعریف مسئله.....	6
2-2-1- یادگیری iid.....	7
2-2-2- تعمیم خارج از توزیع.....	7
2-3- یادگیری بازنمایی علی برای مسئله ی تعمیم خارج از توزیع.....	8
2-3-1- متد.....	10
2-3-1-1- تعمیم دامنه (تعمیم خارج از توزیع) از نگاه علی.....	10
2-3-1-2- یادگیری بازنمایی با الهام از علیت.....	13
2-4- نتیجه گیری.....	18
فصل 3: مروری بر کارهای مرتبط	19
3-1- مقدمه.....	20
3-2- یادگیری بازنمایی بدون نظارت	21
3-2-1- یادگیری بازنمایی تفکیک شده.....	21
3-2-2- یادگیری بازنمایی علی.....	22
3-3- یادگیری مدل با نظارت برای تعمیم خارج از توزیع.....	22
3-3-1- تعمیم دامنه.....	23
3-3-1-1- یادگیری بازنمایی برای تعمیم دامنه.....	23
3-3-1-2- استراتژی یادگیری.....	25
3-3-1-3- تقویت داده.....	26
3-3-2- یادگیری ثابت و علی.....	27
3-3-2-1- روشهای مبتنی بر استنتاج علی.....	27
3-3-2-2- یادگیری ثابت.....	28
3-3-3- یادگیری پایدار.....	29
3-4- بهینه سازی برای تعمیم خارج از توزیع.....	31
3-4-1- بهینه سازی مقاوم توزیعی.....	32
3-4-2- بهینه سازی مبتنی بر عدم تغییر.....	32
3-5- مجموعه داده ها و معیارهای ارزیابی.....	33

33.....	3-5-1- مجموعه داده ها.....
34.....	3-5-1-1- داده های ساختگی.....
35.....	3-5-1-2- داده های واقعی.....
37.....	3-5-2- معیار های ارزیابی.....
39.....	3-6- نتیجه گیری.....

40 فصل 4: نتیجه گیری و کارهای آینده

41.....	4-1- نیاز به محیط های متعدد.....
41.....	4-2- ارزیابی های منطقی.....
41.....	4-3- معرفی موضوع مورد نظر برای پایان نامه.....

42 مراجع

44 واژه نامه

فهرست شکل ها

<u>صفحه</u>	<u>عنوان</u>
9	شکل (1) SCM of DG
13	شکل (2) The framework of CIRL

فهرست جدول‌ها

<u>صفحه</u>	<u>عنوان</u>
36.....	جدول (1) مجموعه داده های تصویری رایج برای تعمیم خارج از توزیع.....

فصل 1:

مقدمه

1-1- شرح مسأله

روش های یادگیری ماشین در بسیاری از زمینه ها مانند پردازش زبان طبیعی، بینایی کامپیوتر و ...، توانایی های زیاد خود را نشان داده اند. از طرفی، بسیاری از تحقیقات جدید، آسیب پذیری مدل های یادگیری ماشین را در مواجهه با داده ها با توزیع های مختلف، نشان داده اند. چنین مشکلی ناشی از نقض یک فرض اساسی است این فرض که داده های آموزش و آزمون به طور یکسان و مستقل توزیع شده اند، که بر اساس آن اکثر مدل های یادگیری موجود توسعه یافته اند. در بسیاری از موارد واقعی جایی که فرض iid^1 را به سختی می توان برآورده کرد، به ویژه آن دسته از کاربرد های پرمخاطره مانند مراقبت های بهداشتی، نظامی و رانندگی خودکار، به جای تعمیم در توزیع آموزشی، توانایی تعمیم تحت تغییر توزیع اهمیت حیاتی تری دارد. بنابراین، بررسی تعمیم خارج از توزیع در هر دو حوزه دانشگاهی و صنعتی از اهمیت زیادی برخوردار است.

علیرغم اهمیت مسئله تعمیم خارج از توزیع، روش های کلاسیک یادگیری با نظارت نمی توانند مستقیماً برای مقابله با مسئله تعمیم خارج از توزیع به کار گرفته شوند. از لحاظ نظری، یکی از اساسی ترین مفروضات یادگیری کلاسیک تحت نظارت، فرض iid^1 است که فرض می کند داده های آموزش و آزمون مستقل و به طور یکسان توزیع شده اند. از آنجایی که، تغییرات توزیعی در مسئله تعمیم خارج از توزیع اجتناب ناپذیر است، باعث می شود فرض iid^1 ارضا نشود و در نتیجه نظریه کلاسیک یادگیری غیرقابل اجر باشد. از نظر تجربی، روش های کلاسیک یادگیری با نظارت معمولاً با کمینه کردن خطا های آموزشی بهینه می شوند، که به صورت حریصانه تمام همبستگی های موجود در داده ها را برای پیش بینی بدست می آورند. اگرچه که در حالت iid^1 موثر است، اما زمانی که تغییرات توزیعی داریم باعث کاهش دقت می شود. زیرا همه همبستگی ها در توزیع های آزمایشی دیده نشده باقی نمی ماند. همانطور که در بسیاری از تحقیقات نشان داده شده است، زمانی که تغییرات توزیعی شدیدی داریم، مدل هایی که صرفاً با خطا های آموزشی بهینه شده اند، به طور چشمگیری عملکرد آن ها کاهش می یابد و گاهی حتی بدتر از حدس های تصادفی عمل می کنند. که این مسائل، ضرورت طراحی روش هایی برای مسئله تعمیم خارج از توزیع را نشان می دهند.

برای مقابله با مسئله تعمیم خارج از توزیع، چندین مسئله ی مهم وجود دارد که باید حل شوند:

- 1) از آنجایی که داده های آموزشی و آزمایشی را می توان از توزیع های مختلف بدست آورد، نحوه توصیف رسمی تغییرات توزیعی هنوز یک مسئله باز است. در مقالات موجود که به موضوع

¹ independent and identically distributed

تعمیم خارج از توزیع پرداخته اند، راه های مختلفی را برای مدل سازی توزیع آزمون بالقوه اتخاذ می کنند. روش های تعمیم دامنه، عمدتاً بر سناریو های واقعی تمرکز می کنند و از داده های جمع آوری شده از حوزه های مختلف استفاده می کنند. روش های یادگیری علی، توزیع های آموزش و آزمون را با ساختارهای علی بیان می کنند و تغییرات توزیعی عمدتاً از تغییرات¹ یا عوامل مخدوش کننده² سرچشمه می گیرند. روش های یادگیری پایدار³ تغییرات توزیعی را از طریق سو گیری انتخاب⁴ معرفی می کنند.

(2) نحوه طراحی یک الگوریتم که عملکرد آن برای مسئله تعمیم خارج از توزیع، مناسب باشد توجه های زیادی را به خود جلب کرده و تحقیقات زیادی روی آن انجام می شود. شاخه های زیادی از روش ها با تمرکز های تحقیقاتی مختلف، از جمله روش های یادگیری بازنمایی بدون نظارت، مدل های یادگیری با نظارت و روش های بهینه سازی وجود دارد.

(3) ارزیابی عملکرد تعمیم خارج از توزیع روش های مختلف همچنان چالش برانگیز است و به مجموعه داده های خاصی با معیارهای ارزیابی مشخصی نیاز دارد، زیرا روش های کلاسیک که برای محیط های iid^5 هستند زمانی که تحت تغییرات توزیعی قرار می گیرند غیرقابل اجرا می شوند. بنابراین این موضوع، انگیزه تولید مجموعه داده ها و ارزیابی های مختلف را ایجاد می کند.

1-2- معرفی حوزه سمینار

در این گزارش، هدف ما ارائه یک بررسی سیستماتیک و جامع از تلاش های تحقیقاتی که در زمینه ی تعمیم خارج از توزیع انجام شده است می باشد که کل چرخه حیات مسئله تعمیم خارج از توزیع را پوشش می دهد. در ابتدا تعریف رسمی مشکل تعمیم خارج از توزیع را ارائه می کنیم. در مرحله بعد از میان روش های موجود، یادگیری علی برای تعمیم دامنه را بررسی می کنیم و یکی از الگوریتم هایی که در این زمینه ارائه شده است به نام الگوریتم CIRL را مرور می کنیم. سپس روش های موجود دیگر را بر اساس موقعیت

¹ Interventions

² confounders

³ Stable Learning

⁴ selection bias

⁵ independent and identically distributed

هایشان در کل مسیر یادگیری به سه بخش، یعنی، یادگیری بازنمایی بدون نظارت، یادگیری مدل تحت نظارت، و بهینه سازی، طبقه بندی می کنیم و روش های رایج برای هر دسته را به تفصیل مورد بحث قرار می دهیم. در نهایت مجموعه داده های رایج و معیارهای ارزیابی را معرفی می کنیم.

3-1- ساختار گزارش

در فصل 2، تعریف مسئله تعمیم خارج از توزیع را ارائه می کنیم و یکی از روش های موجود به نام CIRL که در شاخه ی یادگیری علی برای تعمیم خارج از توزیع قرار دارد را بیان می کنیم. در فصل 3، روش های موجود دیگر را بر اساس موقعیت هایشان در کل مسیر یادگیری در سه بخش ارائه می کنیم و مجموعه داده های رایج و معیارهای ارزیابی که برای مسئله ی تعمیم خارج از توزیع استفاده می شوند را معرفی می کنیم. و در فصل آخر نتیجه گیری نهایی از نتایج این تحقیق ارائه شده است، و برخی مسیر های مناسب برای تحقیق و توسعه بیشتر در این زمینه را پیشنهاد می دهیم.

فصل 2:

تعاریف و مفاهیم مبنایی

1-2- مقدمه

در سال های اخیر با پیچیده شدن مسائل یادگیری ماشین در دنیای واقعی، مسئله ی تعمیم خارج از توزیع، مشکلاتی را برای شبکه های عصبی عمیق مبتنی بر فرض iid^1 به وجود آورده است. به کارگیری مستقیم مدلی که روی دامنه ی مبدا آموزش دیده است به یک دامنه هدف دیده نشده با توزیع متفاوت، باعث کاهش عملکرد الگوریتم می شود که باید به آن توجه شود. برای مقابله با این موضوع مسئله ی تعمیم خارج از توزیع مطرح می شود که یکی از روش های آن، یادگیری بازنمایی علی است.

به جز یادگیری بازنمایی علی برخی از کارها از جمله فرا-یادگیری²، یادگیری بازنمایی ثابت³ و تقویت دامنه⁴ انجام شده است. اگرچه نتایج آن ها قابل قبول است، اما یک مشکل ذاتی در آن ها وجود دارد. این روش ها صرفاً سعی در جبران مشکلات ناشی از داده های خارج از توزیع و مدل سازی وابستگی آماری بین داده ها و برچسب ها بدون توضیح مکانیسم های علی ذاتی دارند. در مقاله ی [1] بحث شده است که چنین اقداماتی ممکن است کافی نباشد، و تعمیم مناسب در خارج از فرض iid' مستلزم یادگیری وابستگی آماری صرف بین متغیر ها نیست، بلکه یک مدل علی ذاتی نیاز است. به عنوان مثال، در یک مسئله طبقه بندی تصویر، به احتمال زیاد همه زرافه ها روی چمن هستند و وابستگی آماری بالایی را نشان می دهند که می تواند به راحتی مدل را گمراه کند و زمانی که پس زمینه در دامنه هدف متفاوت است، پیش بینی اشتباهی انجام دهد. از این گذشته، ویژگی های زرافه ها مانند سر، گردن و ... به جای پس زمینه، یک زرافه را «زرافه» می کند که باید به درستی یاد گرفته شوند و تعیین کننده برچسب هدف باشند.

2-2- تعریف مسئله

فرض کنید که \mathcal{X} فضای ویژگی ها، و \mathcal{Y} فضای برچسب ها باشد. یک مدل دارای پارامتر⁵ به صورت $f_{\theta}: \mathcal{X} \rightarrow \mathcal{Y}$ تعریف می شود که به عنوان یک تابع نگاشت از ویژگی های اصلی به برچسب با پارامتر قابل یادگیری θ عمل می کند. یک تابع زیان به صورت $l: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ تعریف می شود که فاصله بین برچسب پیش بینی شده و برچسب واقعی⁶ را اندازه می گیرد. حال با توجه به این توضیحات می توانیم مسئله کلاسیک یادگیری با نظارت را تعریف کنیم.

¹ independent and identically distributed

² Meta-Learning

³ Invariant Representation Learning

⁴ Domain Augmentation

⁵ parametric

⁶ ground-truth

مسئله 1 (یادگیری با نظارت). یک مجموعه ایی از n نمونه آموزشی به فرم $\{(x_1, y_1), \dots, (x_n, y_n)\}$ داده شده است که دارای توزیع آموزشی $P_{tr}(X, Y)$ هستند. یک مسئله ی یادگیری با نظارت این است که یک مدل بهینه ی f_θ^* را به گونه ایی پیدا کنیم که بتواند تعمیم خوبی روی داده های آزمون که دارای توزیع $P_{te}(X, Y)$ هستند، داشته باشد:

$$f_\theta^* = \underset{f_\theta}{\operatorname{argmin}} E_{X,Y \sim P_{te}} [l(f_\theta(X), Y)] \quad (1)$$

1-2-2- یادگیری iid^1

الگوریتم های یادگیری کلاسیک معمولاً فرض می کنند که نمونه های آموزشی و نمونه های آزمون هر دو تحقق iid^1 از یک توزیع پایه مشترک هستند، که به معنای $P_{tr}(X, Y) = P_{te}(X, Y)$ است. بر اساس چنین فرضیه ای، ERM^2 که میانگین زیان نمونه های آموزشی را به حداقل می رساند، می تواند یک مدل بهینه را به دست آورد که با موفقیت به توزیع آزمون تعمیم یابد. در واقع، عبارت زیر را کمینه می کند:

$$l_{ERM} = \frac{1}{n} \sum_{i=1}^n l(f_\theta(x_i), y_i) \quad (2)$$

2-2-2- تعمیم خارج از توزیع

ویژگی خوب ارائه شده توسط فرض iid^1 زمینه مناسبی را برای توسعه مدل های یادگیری زیادی در چند دهه اخیر فراهم کرده است. با این حال، در موارد واقعی، توزیع داده های آزمون ممکن است از توزیع داده های آموزشی متفاوت باشد، یعنی $P_{tr}(X, Y) \neq P_{te}(X, Y)$. تغییر توزیع می تواند به دلایل زیادی مانند تغییرات زمانی/مکانی داده ها یا سوگیری³ انتخاب نمونه در فرآیند جمع آوری داده ها باشد. در هر صورت، مسئله 1 را پیچیده تر از سناریوی یادگیری $i.i.d$ می کند. علاوه بر این، توزیع داده های آزمون که با آن مواجه می شویم ممکن است به دلیل ماهیت برنامه هایی مانند سناریو آنلاین مبتنی بر جریان⁴ که در آن داده های آزمون در آینده تولید می شوند ناشناخته باشد. به طور خلاصه، مسئله تعمیم خارج از توزیع را می توان به عنوان نمونه ای از مسئله یادگیری با نظارت تعریف کرد که

¹ independent and identically distributed

² Empirical Risk Minimization

³ bias

⁴ stream

در آن توزیع آزمون $P_{te}(X, Y)$ نسبت به توزیع آموزشی $P_{tr}(X, Y)$ دچار تغییر¹ می شود و در طول مرحله آموزش ناشناخته باقی می ماند.

به طور کلی، مسئله تعمیم خارج از توزیع امکان پذیر نیست، مگر اینکه برخی از فرضیات را در مورد چگونگی تغییر توزیع داده های آزمون داشته باشیم. در میان تغییرات توزیعی مختلف، تغییر متغیر² متداول ترین نوع آن است و به طور گسترده مورد مطالعه قرار گرفته است. در واقع تغییر متغیر³ به صورت $P_{tr}(Y|X) = P_{te}(Y|X)$ و $P_{tr}(X) \neq P_{te}(X)$ تعریف می شود، که به این معنی است که توزیع حاشیه ای X از فاز آموزش به فاز آزمون تغییر می کند در حالی که مکانیسم تولید برچسب بدون تغییر باقی می ماند. در این سمینار، تمرکز ما بر روی تغییر متغیر³ است و کار های مختلفی را که در این زمینه انجام شده است، بررسی می کنیم.

2-3- یادگیری بازنمایی علی برای مسئله ی تعمیم خارج از توزیع

در این سمینار ما یک مدل علی ساختاری³ را برای مسئله تعمیم خارج از توزیع با هدف بدست آوردن مکانیزم علی ذاتی بین داده ها و برچسب ها و همینطور بدست آوردن توانایی تعمیم بهتر، معرفی می کنیم. اطلاعات مربوط به هر دسته⁴ در داده ها به عنوان عوامل علی⁵ در نظر گرفته می شود که رابطه آنها با برچسب، مستقل از دامنه است، به عنوان مثال، "شکل" در مسئله تشخیص رقم. در حالی که اطلاعات مستقل از دسته⁶ به عنوان عوامل غیر علی⁷ در نظر گرفته می شود، که عموماً اطلاعات مربوط به دامنه است، به عنوان مثال، "سبک دست خط" در تشخیص رقم. هر داده خام X از ترکیبی از عوامل علی S و عوامل غیرعلی U ساخته شده است و فقط S ها به طور علی بر برچسب دسته Y تأثیر می گذارد، همانطور که در شکل 1 نشان داده شده است. هدف استخراج عوامل علی S از ورودی خام X و سپس بازسازی مکانیزم های علی ثابت⁸ است، که می توان با کمک ایجاد تغییرات علی⁹ $P(Y|do(U), S)$ ، این کار را انجام داد. $do(\cdot)$ نشان دهنده ایجاد تغییرات بر روی متغیر ها است.

¹ shift

² covariate shift

³ SCM; structural causal model

⁴ category-related information

⁵ causal factors

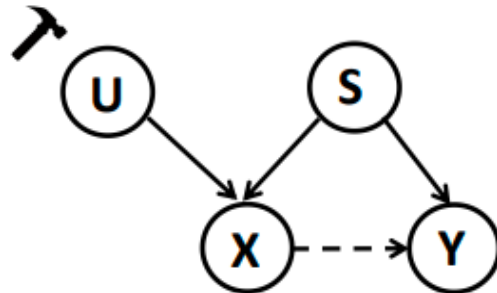
⁶ category

⁷ non-causal factors

⁸ invariant causal mechanisms

⁹ causal intervention

متأسفانه، ما نمی توانیم مستقیماً ورودی خام را به عنوان $X = f(S, U)$ بیان کنیم، زیرا عوامل علی/غیر علی معمولاً مشاهده نمی شوند و نمی توان آنها را فرمول بندی کرد که این موضوع استنتاج علی را چالش برانگیز می کند.



شکل 1: SCM برای *Domain Generalization (Out-Of-Distribution Generalization)* پیکان های توپر نشان می دهد که گره والد باعث و عامل یک فرزند است در حالی که پیکان های خط تیره به معنای وجود وابستگی آماری است.

عوامل علی S انتظار می رود بر اساس [1], [2], [3] سه ویژگی را برآورده کنند:

- (1) جدا و مجزا از عوامل غیر علی U باشند.
- (2) فاکتورسازی S باید به صورت توام مستقل باشد.
- (3) از نظر علی برای مسئله طبقه بندی $X \rightarrow Y$ حاوی تمام اطلاعات علی کافی باشد.

ترکیب با U باعث می شود که S حاوی اطلاعات غیر علی اساسی باشد، در حالی که فاکتورگیری وابسته توام، S را زائد می کند و منجر به از دست رفتن برخی از اطلاعات علی اساسی می شود. در مقابل، عوامل علی S ، عواملی ایده آل هستند که همه الزامات را برآورده می کنند. با الهام از این، یک الگوریتم یادگیری بازنمایی الهام گرفته از علیت¹ پیشنهاد شده است که بازنمایی های آموخته شده را وادار می کند تا ویژگی های بالا را داشته باشند و سپس از هر بعد از بازنمایی ها برای تقلید از فاکتورسازی عوامل علی که توانایی تعمیم قوی تری دارند، بهره برداری می کند.

به طور خلاصه، برای هر ورودی، ابتدا از یک ماژول ایجاد تغییر علی² برای جداسازی عوامل علی S از عوامل

¹ CIRL; Causality Inspired Representation Learning
² causal intervention module

غیرعلی U با تولید داده های جدید با اطلاعات مرتبط با دامنه ی دچار تغییر شده¹ استفاده می کنند. داده های تولید شده دارای عوامل غیر علی متفاوت U هستند اما عوامل علی S در مقایسه با موارد اصلی یکسان هستند، بنابراین بازنمایی ها مجبور می شوند ثابت بمانند. علاوه بر این، آن ها یک مازول فاکتورسازی را پیشنهاد می کنند که هر بُعد از بازنمایی ها را به طور توأم مستقل می کند و سپس می توان از آن برای تقریب عوامل علی استفاده کرد. علاوه بر این، برای اینکه الگوریتم پیشنهادی از نظر علی برای طبقه بندی کافی باشد، یک مازول ماسک متخاصم² طراحی می کنند که به طور مکرر ابعادی را که حاوی اطلاعات علی نسبتاً کمتری هستند شناسایی می کند و آنها را مجبور می کند از طریق یادگیری خصمانه بین یک پوشش دهنده³ و مولد بازنمایی، حاوی اطلاعات علی بیشتر و جدید باشند.

1-3-2- متد

در این بخش، مطابق شکل 1، تعمیم دامنه را از نمای علی با یک مدل علی ساختاری کلی در نظر می گیرند و نشان می دهند که مکانیسم های علی ذاتی (که به عنوان توزیع های شرطی بیان می شوند) می توانند در صورت داشتن عوامل علی، امکان پذیر باشند. با این حال، همانطور که در قسمت قبل گفته شد، بازیابی دقیق عوامل علی دشوار است زیرا آنها غیر قابل مشاهده هستند. بنابراین، پیشنهاد می شود که بازنمایی های علی را بر اساس ویژگی های عوامل علی به عنوان یک تقلید، در حالی که توانایی تعمیم بهتری را دارا می باشند، یاد بگیریم.

1-3-2-1- تعمیم دامنه (تعمیم خارج از توزیع) از نگاه علی

کارهایی که برای تعمیم دامنه انجام شده است بر مدل سازی وابستگی آماری بین ورودی های مشاهده شده و برجسب های مربوطه، یعنی $P(X, Y)$ متمرکز است که در دامنه ها مختلف، متفاوت در نظر گرفته می شود. برای به دست آوردن یک وابستگی ثابت، عموماً توزیع را وادار می کنند که به صورت حاشیه ای یا شرطی ثابت-دامنه⁴ باشد. به عنوان مثال، فاصله بین دامنه ها را در $P(X)$ یا $P(X / Y)$ کمینه می کنند. با این حال، از آنجایی که وابستگی آماری نمی تواند مکانیسم علی ذاتی بین ورودی ها و برجسب ها را توضیح دهد، تمایل به تغییر با دامنه را دارد. بنابراین، وابستگی ثابت یاد گرفته شده در میان دامنه

¹ *perturbed domain-related information*

² *adversarial mask module*

³ *masker*

⁴ *domain-invariant marginally or conditionally*

های مبدا ممکن است همچنان در یک دامنه هدف دیده نشده به خوبی عمل نکنند. در همین حال، مکانیسم های علی معمولاً در همه ی دامنه ها، پایدار و ثابت می مانند [1]. ابتدا ارتباط بین علیت و وابستگی آماری را همانطور که رایشنباخ [2] در اصل 1 ادعا کرده است، بیان می کنیم.

اصل 1. اصل علت مشترک¹: اگر دو متغیر قابل مشاهده X و Y از نظر آماری وابسته باشند، آنگاه یک متغیر S وجود دارد که به طور علی بر هر دو تأثیر می گذارد و تمام وابستگی آن ها را با مستقل کردن آنها تحت شرط S ، بیان می کند.

بر اساس اصل 1 آن ها مدل علی ساختاری زیر را برای توصیف مسئله تعمیم دامنه، بیان می کنند:

$$\begin{aligned} X &:= f(S, U, V_1), S \perp U \perp V_1, \\ Y &:= h(S, V_2) = h(g(x), V_2), V_1 \perp V_2. \end{aligned} \quad (10)$$

که X و Y به ترتیب نشان دهنده ی تصاویر ورودی و برچسب مربوطه ی آن ها است. S نشان دهنده ی عامل های علی است که به صورت علی X و Y را یعنی اطلاعات مربوط به دسته² مانند "شکل" در مسئله ی تشخیص عدد است. درحالی که U نشان دهنده ی فاکتور های غیر علی است که به صورت علی فقط روی X اثر می گذارند. فاکتور های غیر علی در واقع اطلاعات مرتبط با دامنه، مثل "سبک نوشتن" می باشند. V_1 و V_2 ، متغیر های نویز هستند که به صورت توأم مستقل می باشند. f ، h و g را می توان به عنوان توابع ساختاری ناشناخته در نظر گرفت. بنابراین، برای هر توزیع $P(X, Y) \in \mathcal{P}$ ، اگر فاکتور های علی S داده شده باشند یک توزیع شرطی کلی $P(Y|S)$ یعنی یک مکانیزم علی ثابت وجود دارد. طبق صحبت فوق، اگر ما بتوانیم به فاکتور های علی دسترسی پیدا کنیم آنگاه بدست آوردن مکانیزم علی ایی که بتواند به خوبی خارج از فرض iid^3 با بهینه سازی h تعمیم یابد آسان است:

$$h^* = \underset{h}{\operatorname{argmin}} E_P[l(h(g(X)), Y)] = \underset{h}{\operatorname{argmin}} E_P[l(h(S), Y)], \quad (11)$$

که $l(.,.)$ تابع زیان آنتروپی متقابل⁴ است.

مشکلی که وجود دارد این است که ما فاکتور های علی S را از قبل نداریم و تنها تصاویر خام X را داریم که

¹ Common Cause Principle
² category-related information
³ independent and identically distributed
⁴ cross-entropy loss

به صورت کلی بدون ساختار هستند. بازسازی مستقیم عوامل علی و نیز مکانیسم ها به نوعی غیر عملی است، زیرا آنها غیرقابل مشاهده هستند و به صورت مشخص تعریف نشده اند. آنچه که واضح است این است که فاکتور های علی باید از الزامات خاصی پیروی کنند. [1، 3] بیان می کنند که فاکتور های علی باید به صورت توام مستقل باشند، همانطور که در اصل 2 بیان شده است.

اصل 2. اصل مکانیزم های علی مستقل¹:

توزیع شرطی هر متغیر به شرط علل آن (یعنی مکانیسم آن) مکانیسم های دیگر را تحت تأثیر قرار نمی دهد.

از آنجایی که S در معادله (10) نشان دهنده مجموعه همه فاکتور های علی $\{S_1, S_2, \dots, S_N\}$ است، این اصل به ما می گوید که:

- 1) ایجاد تغییر روی یک مکانیسم $P(S_i|PA_i)$ هیچ یک از مکانیسم های دیگر $P(S_j|PA_j)$ را تغییر نمی دهد $i \neq j$ نشان دهنده ی پدر S_i در گراف علی است) که می تواند به عنوان اطلاعات علی در نظر گرفته شود که S_i حاوی است زیرا S پیش از این، گره ریشه است).
- 2) دانستن برخی مکانیسم های دیگر $P(S_i|PA_i)$ اطلاعاتی در مورد مکانیسم $P(S_j|PA_j)$ به ما نمی دهد. بنابراین، می توان توزیع توام فاکتور های علی را به صورت شرطی بیان کرد که به فاکتور سازی علی² اشاره دارد:

$$P(S_1, S_2, \dots, S_N) = \sum_{i=1}^N P(S_i|PA_i), \quad (12)$$

بنابراین خاطر نشان می شویم که فاکتور های علی S باید بر اساس اصل 1 و 2، سه ویژگی پایه را دارا باشند:

- فاکتور های علی S باید جدا و مجزا از فاکتور های غیر علی U باشند یعنی $S \perp U$. بنابراین انجام یک تغییر بر روی U نباید باعث تغییر در S شود.
- فاکتور سازی S_1, S_2, \dots, S_N باید به صورت توام مستقل باشد و هیچکدام از آن ها حاوی اطلاعاتی از بقیه فاکتور ها نباشد.
- فاکتور های علی S باید از نظر علی برای مسئله ی طبقه بندی $X \rightarrow Y$ کافی باشند. یعنی حاوی

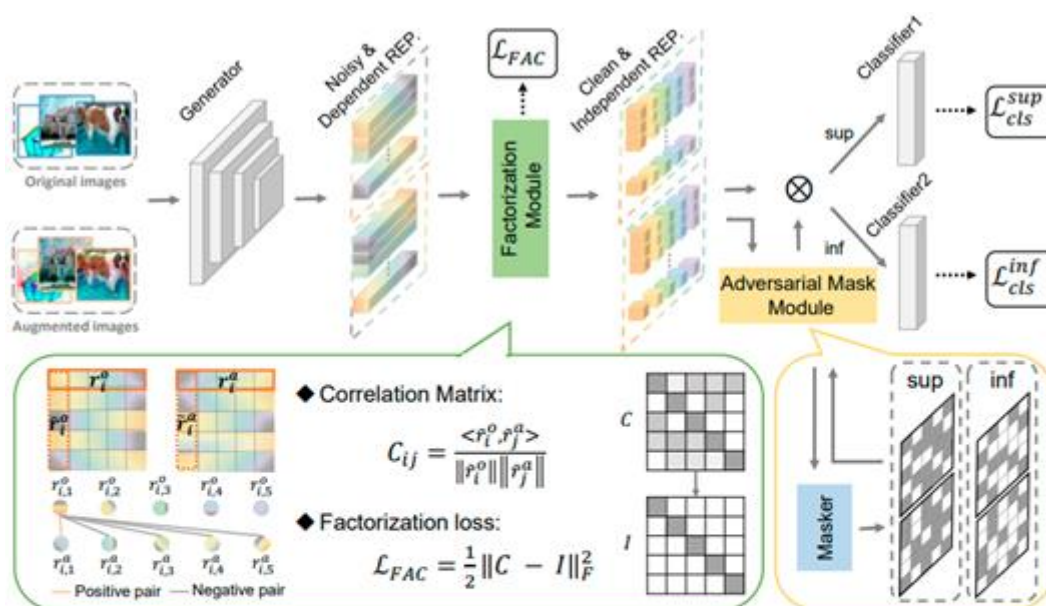
¹ ICM; Independent Causal Mechanisms
² causal factorization

اطلاعاتی باشند که بتواند تمام وابستگی های آماری را توضیح دهد.

بنابراین، به جای بازسازی مستقیم عوامل علی، آن ها پیشنهاد به یادگیری بازنمایی های علی به عنوان جایگزینی با وادار کردن آنها به داشتن ویژگی های یکسان با عوامل علی، می کنند. جزئیات این موضوع را در بخش بعد توضیح خواهیم داد.

2-3-1-2- یادگیری بازنمایی با الهام از علیت¹

در این بخش، الگوریتم یادگیری بازنمایی پیشنهادی الهام گرفته شده از علیت [4] مورد بحث در بالا را نشان می دهیم که از سه ماژول تشکیل شده است: ماژول ایجاد تغییرات علی²، ماژول فاکتور سازی علی³، و ماژول ماسک متخاصم⁴. کل الگوریتم در شکل 3 نشان داده شده است.



شکل 2: چارچوب کلی الگوریتم CIRL. ابتدا تصاویر تقویت شده (*augmented images*) را به وسیله ماژول ایجاد تغییرات علی با ایجاد تغییرات بر عوامل غیر علی تولید می کند. هر دو بازنمایی تصاویر اصلی و تقویت شده به ماژول فاکتور سازی ارسال می شوند، که یک تابع زیان فاکتور سازی را تحمیل می کند تا بازنمایی ها از فاکتور های غیر علی جدا و مجزا و همچنین به طور توأم مستقل شوند. در نهایت، ماژول ماسک متخاصم، رقابت را بین مولد و یک پوشش دهنده (*masker*) ایجاد می کند و بازنمایی ها را به طور علی کافی برای طبقه بندی ارائه می دهد.

¹ Causality Inspired Representation Learning

² causal intervention module

³ causal factorization module

⁴ adversarial mask module

ماژول ایجاد تغییرات علی¹

در ابتدا هدف جدا سازی فاکتور های علی S از ترکیب آن ها با فاکتور های غیر علی U به وسیله ی ایجاد تغییرات علی است. اگرچه فرم صریح استخراج کننده فاکتور های علی $g(.)$ در عبارت (11) به صورت کلی ناشناخته است اما با توجه به دانش قبلی می دانیم که فاکتور های علی S باید هنگام ایجاد تغییرات بر روی U ، ثابت بمانند (یعنی $(P(S | do(U)))$ ، در حالی که از تحقیقات مربوط به تعمیم دامنه (تعمیم خارج از توزیع) می دانیم که بعضی از اطلاعات مرتبط با دامنه نمی توانند دسته تصویر ورودی را مشخص کنند که به عنوان فاکتور های غیر علی در نظر گرفته می شوند و توسط بعضی از تکنیک ها تسخیر می شوند. به طور مثال تبدیل فوریه یک ویژگی مناسب در این زمینه دارد: مولفه فاز طیف فوریه معنای سطح بالای سیگنال اصلی را حفظ می کند، در حالی که مولفه دامنه حاوی اطلاعات آماری سطح پایین است. بنابراین، روش $CIRL$ ایجاد تغییرات را بر روی U با برهم زدن اطلاعات دامنه انجام می دهد در حالی که اطلاعات فاز را بدون تغییر نگه می دارد. با توجه به یک تصویر ورودی اصلی x^0 ، تبدیل فوریه آن را می توان به صورت زیر بیان کرد:

$$\mathcal{F}(x^0) = \mathcal{A}(x^0) \times e^{-j \times \mathcal{P}(x^0)}, \quad (13)$$

که $\mathcal{A}(x^0)$ و $\mathcal{P}(x^0)$ به ترتیب نشان دهنده ی مولفه ی دامنه و فاز است. تبدیل فوریه ی $\mathcal{F}(.)$ و معکوس آن $\mathcal{F}^{-1}(.)$ می تواند به صورت موثر توسط الگوریتم FFT محاسبه شوند. سپس اطلاعات دامنه را از طریق درون یابی خطی بین طیف های دامنه تصویر اصلی x^0 و تصویر $(x')^0$ که به طور تصادفی از دامنه های مبدا دلخواه نمونه برداری شده است، تغییر می دهند:

$$\hat{\mathcal{A}}(x^0) = (1 - \lambda)\mathcal{A}(x^0) + \lambda\mathcal{A}((x')^0), \quad (14)$$

که $\lambda \sim U(0, \eta)$ و η میزان تغییرات را کنترل می کند. سپس طیف های دامنه تغییر یافته را با مولفه فاز اصلی ترکیب می کنند تا تصویر تقویت شده x^a را با تبدیل فوریه معکوس ایجاد کنند:

$$\mathcal{F}(x^a) = \hat{\mathcal{A}}(x^0) \times e^{-j \times \mathcal{P}(x^0)}, \quad x^a = \mathcal{F}^{-1}(\mathcal{F}(x^a)). \quad (15)$$

مولد بازنمایی پیاده سازی شده توسط یک مدل CNN را با $\hat{g}(\cdot)$ و بازنمایی ها را با $r = \hat{g}(x) \in \mathbb{R}^{1 \times N}$ نشان می دهیم که در آن N تعداد ابعاد است. برای شبیه سازی فاکتور های علی که با ایجاد تغییرات در U ثابت می مانند، \hat{g} را بهینه سازی می کنیم تا بازنمایی ها را در هر بُعد با ایجاد تغییرات بالا بدون تغییر نگه دارند:

$$\max_{\hat{g}} \frac{1}{N} \sum_{i=1}^N COR(\tilde{r}_i^o, \tilde{r}_i^a), \quad (16)$$

که \tilde{r}_i^o و \tilde{r}_i^a به ترتیب نشان دهنده ی Z -score نرمال شده i امین ستون $R^o = [(r_1^o)^T, \dots, (r_B^o)^T]^T$ و $R^a = [(r_1^a)^T, \dots, (r_B^a)^T]^T \in \mathbb{R}^{B \times N}$ است. $B \in \mathbb{Z}_+$ اندازه دسته ها¹، $r_i^o = \hat{g}(x_i^o)$ و $r_i^a = \hat{g}(x_i^a)$ برای $i \in \{1, \dots, B\}$ است. از یک تابع COR برای اندازه گیری همبستگی بازنمایی ها قبل و بعد از ایجاد تغییر استفاده می کنند. بنابراین، می توان اولین مرحله شبیه سازی فاکتور های علی S با بازنمایی های R را با مستقل کردن آنها از U درک کرد.

ماژول فاکتور سازی علی²

همانطور که در بخش "تعمیم دامنه از نگاه علی" گفته شد که فاکتور سازی فاکتور های علی S_1, S_2, \dots, S_N باید به صورت توام مستقل باشند طوری که هیچ یک از آنها حاوی اطلاعات دیگران نیست. بنابراین، در الگوریتم $CIRL$ تلاش می کنند که هر دو بعد از بازنمایی ها را مستقل از یکدیگر کنند:

$$\min_{\hat{g}} \frac{1}{N(N-1)} \sum_{i \neq j} COR(\tilde{r}_i^o, \tilde{r}_j^a), i \neq j, \quad (17)$$

برای کم کردن هزینه محاسبات محدودیت های درون R^a یا R^o را در نظر نگرفته اند. برای تحقق اهداف بهینه سازی عبارات (16) و (17)، یک ماتریس همبستگی C به صورت زیر در نظر می گیرند:

¹ batch size
² causal factorization module

$$C_{ij} = \frac{\langle \tilde{r}_i^o, \tilde{r}_j^a \rangle}{\|\tilde{r}_i^o\| \|\tilde{r}_j^a\|}, i, j \in 1, 2, \dots, N, \quad (18)$$

که $\langle \cdot \rangle$ نشان دهنده ی عملگر ضرب داخلی است. بنابراین بُعد های یکسان R^o و R^a می توانند به عنوان جفت نمونه های مثبت در نظر گرفته شوند که باید همبستگی را به حداکثر برسانند. در حالی که بُعد های متفاوت می توانند به عنوان جفت نمونه های منفی در نظر گرفته شوند که باید همبستگی را به حداقل برسانند. بر این اساس، آن ها یک تابع زیان \mathcal{L}_{Fac} در نظر می گیرند که به صورت زیر بیان می شود:

$$\mathcal{L}_{Fac} = \frac{1}{2} \|C - I\|_F^2. \quad (19)$$

معادله (19) می تواند عناصر قطری ماتریس همبستگی C را به صورت تقریبی 1 کند، به این معنی که بازنمایی های قبل و بعد از ایجاد تغییر نسبت به فاکتور های غیر علی، ثابت هستند. این نشان می دهد که می توانیم به طور مؤثر فاکتور های علی را از ترکیب آن ها با فاکتور های غیر علی جدا کنیم. علاوه بر این، عناصر غیر قطری C را به 0 نزدیک می کند، یعنی ابعاد بازنمایی ها را به صورت توام مستقل می کند. بنابراین، با به حداقل رساندن \mathcal{L}_{Fac} ، می توانیم بازنمایی های دارای نویز و وابسته را به بازنمایی هایی تمیز و مستقل تبدیل کنیم که دو ویژگی اول فاکتور های علی ایده آل را برآورده می کنند.

ماژول ماسک متخاصم¹

برای موفقیت در مسئله ی طبقه بندی $X \rightarrow Y$ ، بازنمایی ها باید به صورت علی مناسب باشند به طوری که شامل همه ی اطلاعات کمک کننده باشند. سراسر ترین راه بکار گرفتن برچسب های با نظارت y در چند دامنه مبدا است:

$$\mathcal{L}_{cls} = l(\hat{h}(\hat{g}(x^o)), y) + l(\hat{h}(\hat{g}(x^a)), y) \quad (20)$$

که \hat{h} طبقه بندی کننده است. با این حال، این راه ساده نمی تواند تضمین کند که هر بُعد از بازنمایی های

¹ adversarial mask module

یاد گرفته شده ما مهم است، یعنی حاوی اطلاعات علی اساسی کافی برای طبقه بندی است. ممکن است ابعادی با ارزش پایین تری¹ وجود داشته باشد که اطلاعات علی نسبتاً کمتری را دارا می باشند و در نتیجه سهم کوچکی در طبقه بندی دارند. بنابراین، نویسندگان روش *CIRL* پیشنهاد می کنند که این ابعاد را شناسایی کرده و آنها را برای مشارکت بیشتر به اجرا درآورد. از آنجایی که ابعاد نیز باید به کمک مازول فاکتور سازی ارائه شده به صورت توام مستقل باشند، ابعاد با ارزش پایین تر² شناسایی شده به گونه ای ارائه می شوند که حاوی اطلاعات علی بیشتر و جدید باشند که در سایر ابعاد گنجانده نشده است، که باعث می شود کل بازنمایی ها از جهت علی غنی تر باشند.

بنابراین، برای تشخیص ابعاد با ارزش پایین تر، آن ها یک مازول ماسک متخصص طراحی می کنند. یک ماسک مبتنی بر شبکه عصبی می سازند که با \hat{W} نشان داده می شود تا سهم هر بعد را یاد بگیرد و ابعاد مربوط به بزرگترین نسبت $\kappa \in (0,1)$ به عنوان ابعاد با ارزش بالاتر³ در نظر گرفته می شود. در حالی که بقیه به عنوان ابعاد با ارزش پایین تر در نظر گرفته می شوند:

$$m = \text{Gumbel_Softmax}(\hat{w}(r), \kappa N) \in \mathbb{R}^N, \quad (21)$$

که آن ها از ترفند مشتق پذیر رایج *Gumbel_Softmax* برای نمونه برداری از ماسکی با مقادیر κN نزدیک به 1 استفاده می کنند. با ضرب بازنمایی های یاد گرفته شده در ماسک های بدست آمده m و $1-m$ می توان به ترتیب ابعاد با ارزش بالاتر⁴ و پایین تر⁵ بازنمایی ها را بدست آورد. سپس، آنها را به دو طبقه بندی کننده مختلف \hat{h}_1, \hat{h}_2 وارد می کنند. معادله (20) را می توان به صورت زیر بازنویسی کرد:

$$\begin{aligned} \mathcal{L}_{cls}^{sup} &= l(\hat{h}_1(r^o \odot m^o), y) + l(\hat{h}_1(r^a \odot m^a), y), \\ \mathcal{L}_{cls}^{inf} &= l(\hat{h}_2(r^o \odot (1 - m^o)), y) + l(\hat{h}_2(r^a \odot (1 - m^a)), y), \end{aligned} \quad (22)$$

آن ها ماسک کننده⁶ را با کمینه کردن \mathcal{L}_{cls}^{sup} و بیشینه کردن \mathcal{L}_{cls}^{inf} ، بهینه می کنند. درحالی که مولد \hat{g}

¹ inferior
² inferior dimensions
³ superior dimensions
⁴ superior dimensions
⁵ inferior dimensions
⁶ masker

و طبقه بندی کننده های \hat{h}_1, \hat{h}_2 را با کمینه کردن دو تابع زیان با نظارت، بهینه می کنند.

ماژول ماسک متخاصم به طور دقیق می تواند ابعاد با ارزش پایین تر² را شناسایی کند چون:

(1) برای یک \hat{h}_2 بهینه شده با هدف کمینه کردن \mathcal{L}_{cls}^{inf} بر اساس ابعاد ماسک شده¹ موجود، یادگیری m برای انتخاب ابعاد برای به حداکثر رساندن \mathcal{L}_{cls}^{inf} می تواند ابعاد با ارزش پایین تر که مشارکت کمتری دارند را پیدا کند.

(2) مجموعه ابعاد با ارزش بالاتر¹ و پایین تر² مکمل یکدیگر هستند به طوری که اگر یک بُعد به عنوان بُعد با ارزش بالاتر تلقی نشود، به عنوان یک بُعد با ارزش پایین تر تلقی می شود، بنابراین انتخاب ابعاد با ارزش بالاتر به انتخاب ابعاد با ارزش پایین تر کمک می کند.

علاوه بر این، در مقایسه با بهینه سازی معادله (20)، ماژول ماسک متخاصم همراه با ماژول فاکتورسازی علی می تواند به تولید بازنمایی هایی از جهت علی کافی تر کمک کند، زیرا با بهینه سازی \hat{g} برای کمینه کردن \mathcal{L}_{cls}^{inf} و \mathcal{L}_{Fac} ، ابعاد با ارزش پایین تر مجبور می شوند اطلاعات علی بیشتری را حمل کنند و مستقل از ابعاد با ارزش بالاتر موجود، باشند. در نهایت، بازنمایی های آموخته شده با «جایگزینی» مکرر بازنمایی های با ارزش پایین تر به عنوان بازنمایی های با ارزش بالاتر جدید، به حالتی که از جهت علی کافی تر باشند نزدیک می شوند.

هدف کلی بهینه سازی CIRL به صورت زیر خلاصه می شود:

$$\min_{\hat{g}, \hat{h}_1, \hat{h}_2} \mathcal{L}_{cls}^{sup} + \mathcal{L}_{cls}^{inf} + \tau \mathcal{L}_{Fac}, \quad \min_{\hat{w}} \mathcal{L}_{cls}^{sup} - \mathcal{L}_{cls}^{inf}, \quad (23)$$

که در آن τ یک پارامتر برای تنظیم شدت تاثیر \mathcal{L}_{Fac} است. توجه داشته باشید که کل بازنمایی r و طبقه بندی کننده \hat{h}_1 در طول استنتاج² استفاده می شوند.

4-2- نتیجه گیری

در این فصل ابتدا، مسئله تعمیم خارج از توزیع را به صورت کلی بیان کردیم و ارتباط و تفاوت آن با مسئله یادگیری با فرض iid را نشان دادیم. سپس، دو اصل اساسی از یادگیری علی را بیان کردیم و بر اساس آن یک الگوریتم یادگیری بازنمایی الهام گرفته شده از علیت به نام CIRL را معرفی کردیم.

¹ masked
² inference

فصل 3:

مروري بر کارهاي مرتبط

1-3- مقدمه

برای مقابله با چالش های ناشی از تغییرات توزیعی که ناشناخته است، تلاش های زیادی در تعمیم خارج از توزیع انجام شده است. تکنیک های استفاده شده از جمله علیت، یادگیری بازنمایی، تکنیک های مبتنی بر ساختار و مبتنی بر بهینه سازی، بسیار متفاوت از هم هستند. در این سمینار سعی شده است که به تفصیل این تکنیک ها توضیح داده شوند.

به طور کلی، مسئله یادگیری تحت نظارت تعریف شده در معادله 1 را می توان به سه جزء به صورت نسبتاً مستقل تقسیم کرد:

- 1) بازنمایی ویژگی های X (به عنوان مثال $g(X)$).
- 2) تابع نگاشت $f_\theta(X)$ از ویژگی های X (یا $g(X)$) به برچسب Y ، که عموماً به عنوان مدل یا سوگیری استقرایی¹ نیز شناخته می شود.
- 3) فرمول بندی هدف بهینه سازی.

بنابراین، روش های موجود را بر اساس موقعیت آنها در کل مسیر یادگیری به سه بخش دسته بندی می کنیم:

- 1) یادگیری بازنمایی بدون نظارت برای تعمیم خارج از توزیع: شامل یادگیری بازنمایی تفکیک شده² و یادگیری بازنمایی علی می باشد که از تکنیک های یادگیری بازنمایی بدون نظارت (مثل *variational Bayes*) برای جاسازی³ کردن دانش قبلی در فرایند یادگیری، استفاده می کند.
- 2) یادگیری مدل با نظارت برای تعمیم خارج از توزیع: شامل یادگیری علی، یادگیری پایدار⁴ و تعمیم دامنه می باشد که معماری مدل های مختلف و استراتژی های یادگیری را برای دستیابی به تعمیم خارج از توزیع طراحی می کند.
- 3) بهینه سازی برای تعمیم خارج از توزیع: شامل بهینه سازی مقاوم توزیعی⁵ و بهینه سازی مبتنی بر عدم تغییر⁶ که به صورت مستقیم تعمیم خارج از توزیع را فرموله می کند و بهینه

¹ *inductive bias*

² *Disentangled Representation Learning*

³ *embed*

⁴ *stable learning*

⁵ *Distributionally Robust Optimization*

⁶ *Invariance-Based Optimization*

سازی را با تضمین های نظری برای تعمیم خارج از توزیع در نظر می گیرد.

2-3- یادگیری بازنمایی بدون نظارت

در این بخش، روش هایی که بر روی یادگیری بازنمایی بدون نظارت تمرکز دارند را بررسی می کنیم، که می توان آن ها را عمدتاً به دو شاخه تقسیم کرد، یعنی یادگیری بازنمایی تفکیک شده¹ و یادگیری بازنمایی علی. این روش ها از دانش قبلی انسان برای طراحی رویه و مراحل یادگیری بازنمایی استفاده می کنند، که به بازنمایی یادگرفته شده، ویژگی های خاصی می بخشد که به طور بالقوه برای تعمیم خارج از توزیع مفید هستند.

1-2-3- یادگیری بازنمایی تفکیک شده

هدف یادگیری بازنمایی تفکیک شده، یادگیری بازنمایی هایی است که در آن فاکتور های متمایز و دارای اطلاعات مفید از تغییرات در داده ها، از هم تفکیک شده باشند [5]، که به عنوان یک ویژگی از بازنمایی خوب در نظر گرفته می شود و به طور بالقوه از تعمیم خارج از توزیع بهره می برد. کارهایی که برای انجام تفکیک سازی با روش های مبتنی بر VAE^2 انجام شده اند، هم بر تفسیر پذیری و هم بر پراکندگی³ تأکید می کنند، که پراکندگی به معنای تغییرات کوچک توزیع است که معمولاً خود را به صورت پراکنده یا محلی در فاکتور گیری تفکیک شده⁴ نشان می دهند.

علیرغم موفقیت یادگیری بازنمایی تفکیک شده، لوکاتلو و همکاران [5] برخی از مفروضات رایج یادگیری بازنمایی تفکیک شده بدون نظارت را به چالش می کشد (به عنوان مثال، استقلال عوامل پنهان⁵). همچنین این پرسش را مطرح می کند که آیا تفکیک سازی می تواند عملکردهای *downstream task* را بهبود بخشد، و الهام بخش کارهای بعدی برای در نظر گرفتن *downstream task* ها، از جمله عملکرد تعمیم خارج از توزیع می باشد.

با این حال، اینکه آیا بازنمایی تفکیک شده به نفع تعمیم خارج از توزیع است، بحث برانگیز است. محققان برخی آزمایش های برون یابی کمی را انجام دادند و دریافتند که بازنمایی تفکیک شده آموخته شده به داده های دیده نشده تعمیم نمی یابد. در مجموع، مزیت بازنمایی تفکیک شده در وظایف خارج از توزیع

¹ disentangled

² Variational Autoencoders

³ sparsity

⁴ disentangled factorization

⁵ latent factor

هنوز به تحقیق و بحث بیشتری نیاز دارد.

2-2-3- یادگیری بازنمایی علی

مشابه یادگیری بازنمایی تفکیک شده، هدف یادگیری بازنمایی علی، یادگیری متغیرها در گراف علی به یک روش بدون نظارت یا نیمه نظارتی¹ است. علاوه بر این، بازنمایی علی را می توان به عنوان هدف نهایی روش تفکیک شده در نظر گرفت، که تعریف غیررسمی بازنمایی تفکیک شده را از نظر تفسیر پذیری و پراکندگی برآورده می کند. با بازنمایی علی یادگرفته شده، می توان به فرآیند تولید داده های پنهان² پی برد، که می تواند به پایداری³ در برابر تغییرات توزیعی ناشی از مداخلات⁴ کمک کند.

در سناریو های واقعی که مشاهدات به جای داده های ساختار یافته در قالب تصاویر یا جملات شکل گرفته اند، اطلاعات انتزاعی سطح بالا باید از داده های سطح پایین استخراج شود، و در مقالات موجود پیشنهاد می شود که فاکتورسازی علی⁵ را از طریق تفکیک سازی بدست آورد.

3-3- یادگیری مدل با نظارت برای تعمیم خارج از توزیع

جدای از یادگیری بازنمایی ها صرفاً به صورت بدون نظارت، شاخه هایی از کارهای انجام شده وجود دارد که اطلاعات نظارت شده را برای طراحی معماری های مدل مختلف و استراتژی های یادگیری مربوطه در نظر می گیرد. در این بخش، روش هایی را بررسی می کنیم که بر یادگیری مدل انتها به انتها⁶ برای دستیابی به توانایی تعمیم خارج از توزیع تمرکز می کنند، از جمله تعمیم دامنه، یادگیری علی، یادگیری ثابت⁷ و یادگیری پایدار⁸.

¹ semi-supervised

² the latent data generation process

³ resist

⁴ interventions

⁵ causal factorization

⁶ end-to-end

⁷ invariant learning

⁸ stable learning

1-3-3- تعمیم دامنه

هدف تعمیم دامنه¹، یادگیری مدل از طریق ترکیب داده‌ها از چندین دامنه مختلف است که به خوبی در دامنه هدف دیده نشده تعمیم یابد که عمدتاً بر مسائل طبقه‌بندی مربوط به بینایی کامپیوتر تمرکز دارد، زیرا پیش‌بینی‌ها تحت تاثیر بر اختلال² در تصاویر (مانند سبک، پس زمینه، نور، چرخش و غیره) هستند. روش‌های تعمیم دامنه را به سه شاخه، یعنی یادگیری بازنمایی، استراتژی آموزش و تقویت داده‌ها تقسیم می‌کنیم که به اختصار آن‌ها را معرفی خواهیم کرد.

1-3-1- یادگیری بازنمایی برای تعمیم دامنه

یادگیری بازنمایی یک بخش مهم در تعمیم دامنه است. مقاله‌ی [6] از نظر تئوری یا تجربی ثابت می‌کند که اگر بازنمایی‌ها در هنگام تغییر دامنه ثابت بمانند، بازنمایی‌ها قابل انتقال³ و مقاوم⁴ به دامنه‌های مختلف هستند. روش‌هایی که تلاش می‌کنند بازنمایی‌های ثابت را در میان دامنه‌های مختلف یاد بگیرند، عمدتاً به سه دسته تقسیم می‌شوند:

1. یادگیری متخاصم دامنه⁵

2. هم ترازی دامنه⁶

3. روش‌های مبتنی بر هسته.

یادگیری متخاصم دامنه: محققان یک شبکه عصبی متخاصم دامنه⁷ را برای تطبیق دامنه⁸ پیشنهاد کرده‌اند. *DANN* با بهینه‌سازی توأم ویژگی‌های پایه، بازنمایی‌هایی را می‌آموزد که با تغییر دامنه، متمایز و ثابت هستند. همچنین *DANN* دارای یک پیش‌بینی کننده برچسب است که برچسب‌های کلاس را پیش‌بینی می‌کند که هم در مرحله آموزش و هم در مرحله استنتاج⁹ استفاده می‌شود. همینطور *DANN* دارای یک

¹ DG; Domain Generalization

² disturbance

³ transferable

⁴ robust

⁵ domain adversarial learning

⁶ domain alignment

⁷ DANN

⁸ domain adaptation

⁹ inference phase

طبقه بندی دامنه است که بین دامنه های مبدا و هدف در طول آموزش تمایز قائل می شود. در واقع بازنمایی ها به گونه ای آموزش داده می شوند که طبقه بندی کننده دامنه را گیج کنند تا ویژگی های ثابت دامنه یاد گرفته شود. لی و همکاران [7] این ایده را برای تعمیم دامنه استفاده کرده اند. و گنگ و همکاران [7] آموزش متخاصم را بیشتر به فضای چندگانه¹ گسترش داده اند. لی و همکاران [7] پیشنهاد به یادگیری شبکه های متخاصم کلاس خاص² از طریق یک شبکه متخاصم ثابت شرطی³ می کنند.

هم ترازی دامنه:⁴ برخی از مقاله ها پیشنهاد به یادگیری بازنمایی های ثابت دامنه⁵ از طریق هم ترازی⁶ ویژگی ها می کنند. معطیان و همکاران [7] پیشنهاد به یادگیری هم ترازی معنایی⁷ بین دامنه های مختلف با به حداقل رساندن فاصله بین نمونه ها از دامنه های مختلف اما کلاس یکسان و به حداکثر رساندن فاصله بین نمونه ها از دامنه ها و کلاس های مختلف می کنند. برخی از مقاله ها با به حداقل رساندن فاصله حداکثر میانگین اختلاف⁸، فاصله *Wasserstein* و همبستگی مرتبه دوم، واگرایی توزیع ویژگی ها را برای تطبیق دامنه یا تعمیم دامنه به حداقل می رسانند.

روش های مبتنی بر هسته: روش های مبتنی بر هسته نیز به طور گسترده در تعمیم دامنه استفاده می شوند. بلانچارد و همکاران [7] ابتدا مسئله تعمیم دامنه را از طریق روش های مبتنی بر هسته بررسی می کنند و پیشنهاد به یادگیری یک هسته ثابت دامنه⁹ با داده های آموزشی می کنند. مواندت و همکاران [7] یک روش مبتنی بر هسته کلاسیک برای تعمیم دامنه به نام تجزیه و تحلیل مؤلفه های ثابت دامنه¹⁰ پیشنهاد می کنند که واریانس توزیعی بین نمونه ها از دامنه های مبدا را به حداقل می رساند. گان و همکاران [7] *DICA* را با تنظیم صفت¹¹ گسترش می دهند.

¹ manifold space

² class-specific adversarial networks

³ CIAN- conditional invariant adversarial network

⁴ domain alignment

⁵ domain invariant representations

⁶ alignment

⁷ semantic alignment

⁸ MMD-maximum mean discrepancy

⁹ domain-invariant kernel

¹⁰ DICA- Domain-Invariant Component Analysis

¹¹ attribute regularization

2-1-3 استراتژی یادگیری

با توجه به اینکه تعمیم دامنه بیشتر بر مسائل طبقه بندی مربوط به بینایی کامپیوتری تمرکز می کند، برخی از مقالات، استراتژی های یادگیری را برای افزایش توانایی تعمیم مدل های عمیق بر روی داده های تصویری به کار می گیرند که می توان آنها را به چهار دسته تقسیم کرد: فرا یادگیری¹، یادگیری گروهی²، تعمیم دامنه بدون نظارت/نیمه-نظارتی³ و مابقی روش ها.

فرا یادگیری:

فرا یادگیری از طریق یک الگوی یادگیری جایگزین، تجربه و دانش را در طول دوره های یادگیری متعدد به دست می آورد. فین و همکاران [7] فرایادگیری مدل-آگنوستیک⁴ را برای انطباق دامنه⁵ پیشنهاد می کنند، که مفهوم «قسمت ها»⁶ را در مرحله آموزش معرفی می کند و تا حد زیادی بر تحقیقات فرا یادگیری برای تعمیم دامنه تأثیر می گذارد. ایده اصلی فرا یادگیری برای تعمیم دامنه این است که دامنه های مبدا را به فرا آموزش⁷ و فرا آزمون⁸ تقسیم کنیم، جایی که تابع زیان فرا آموزش⁹ و فرا آزمون¹⁰ به طور همزمان بهینه می شوند.

یادگیری گروهی:

به طور معمول، روش های مبتنی بر یادگیری گروهی، مجموعه ای از چندین مدل خاص را برای دامنه های مبدا مختلف یاد می گیرند تا توانایی تعمیم را بهبود بخشند. برخی از مقاله ها، زیر شبکه های خاص-دامنه¹¹ را برای دامنه های مبدا مختلف، همراه با یک طبقه بندی کننده واحد یا چندین سر طبقه بندی کننده خاص-دامنه در نظر می گیرند. برخی دیگر از نرمال سازی دسته ای خاص-دامنه برای دامنه های مختلف استفاده می کنند تا نرمال سازی بهتری را بیاموزند.

¹ meta learning

² ensemble learning

³ unsupervised/semi-supervised DG

⁴ MAML; model-agnostic meta-learning

⁵ Domain Adaptation

⁶ episodes

⁷ meta-train

⁸ meta-test

⁹ meta-train

¹⁰ meta-test

¹¹ domain-specific

تعمیم دامنه بدون نظارت/نیمه-نظارتی:

اخیراً، با الهام از انطباق¹ دامنه بدون نظارت و نیمه نظارت، برخی از مقاله ها پیشنهاد می کنند که توانایی تعمیم مدل را با یادگیری بدون نظارت یا نیمه نظارت افزایش دهند. ژانگ و همکاران [7] آموزش بدون نظارت نامربوط دامنه² را برای مقابله با تغییرات توزیع بین دامنه های مبدا و دامنه های هدف پیشنهاد می کنند. آنها نمونه های منفی معتبر را برای هر نمونه صف داده شده با توجه به شباهت بین دامنه های مختلف انتخاب می کنند تا بازنمایی های نامربوط دامنه³ را یاد بگیرند.

مابقی روش ها:

کارلوجی و همکاران با الهام از روش یادگیری خود نظارتی⁴ [7] یک مسئله پازل خودآموز را با مسئله طبقه بندی ترکیب می کنند تا بازنمایی های مقاوم⁵ را بیاموزند. ریو و همکاران [7] نمونه برداری از نمونه های مثبت و منفی را از طریق جنگل تصادفی⁶ پیشنهاد می کنند. لی و همکاران [7] به طور متناوب لایه های کانولوشن و طبقه بندی کننده را آموزش می دهند. هوانگ و همکاران [7] یک الگوریتم حذف تصادفی خود چالش برانگیز⁷ را برای جلوگیری از بیش برآزش مدل به دامنه های مبدا معرفی می کنند.

3-3-1-3 تقویت داده

تقویت داده ها روشی رایج و موثر در یادگیری عمیق به ویژه در بینایی کامپیوتری است. توانایی تعمیم مدل های عمیق تا حد زیادی به ناهمگونی داده های موجود بستگی دارد. بنابراین ناهمگونی ناشی از افزایش داده ها می تواند از بیش برآزش جلوگیری کند و توانایی تعمیم را بهبود بخشد. روش های افزایش داده ها برای تعمیم دامنه را می توان به تقویت مبتنی بر تصادفی بودن⁸، تقویت مبتنی بر گرادیان⁹ و تقویت مبتنی بر

¹ adaptation

² DIUL; Domain-Irrelevant Unsupervised Learning

³ domain-irrelevant representations

⁴ self-supervised

⁵ robust

⁶ random forest

⁷ self-challenging dropout algorithm

⁸ randomization based augmentation

⁹ gradient-based augmentation

تولید¹ تقسیم کرد [8], [9], [10].

2-3-3- یادگیری ثابت و علی²

در مقایسه با تعمیم دامنه که به طور معمول مسائل بینایی را هدف قرار می دهد ، یادگیری علی و یادگیری ثابت ، از موضوع استنتاج علی ناشی می شود و مسئله تعمیم خارج از توزیع را به روشی اصولی تر بررسی می کند، که هدف آن کشف متغیرهای علی برای پیش بینی است و اخیراً کاربردی تر شده است. برای روش های یادگیری علی، اغلب فرض می شود که ناهمگونی داده و رابطه علی در داخل داده ها وجود دارد. به طور خاص، یادگیری علی فرض می کند که فرد به داده ها از چندین محیط دسترسی دارد.

1-2-3- روشهای مبتنی بر استنتاج علی

ابتدا روش های مربوط به استنتاج علی را بررسی می کنیم که سعی می کند متغیرهای علی را از داده های ناهمگن به دست آورد.

به وضوح مشخص است که یک استاندارد مناسب برای شناسایی اثر علی یک متغیر، انجام آزمایش های تصادفی سازی شده مانند آزمایش A/B است، اما آزمایش های کاملاً تصادفی معمولاً گران هستند و حتی در کاربردهای واقعی غیرممکن هستند. از آنجایی که استنتاج علی³ یا یادگیری ساختاری علی⁴ بسیار بلندپروازانه است، این تکنیک های استنتاج باید به عنوان "حقیقت پایه"⁵ در نظر گرفته شوند، اما لزوماً نمی توانند در عمل محقق شوند. بنابراین، طراحی چنین تکنیک هایی که از نظر "توضیح علی" نسبت به رگرسیون استاندارد⁶ یا چارچوب طبقه بندی⁷ غنی تر هستند و همچنین می توانند نوعی تغییر ناپذیری را در بین محیط ها به دست آورند، عملی تر است. به دنبال چنین شهودی، رشته ای از روش ها [11] با بهره بردن از ناهمگونی درون داده ها (به عنوان مثال، محیط های متعدد) توسعه یافته اند.

¹ generation based augmentation

² Causal & Invariant Learning

³ causal inference

⁴ causal structural learning

⁵ ground truth

⁶ standard regression

⁷ classification framework

پیترز و همکاران [7] در ابتدا سعی کردند که این واقعیت را بررسی کنند که «عدم تغییر»¹ تا حدی می تواند ساختار علی را تحت شرایط لازم استنتاج کند و پیش بینی علی ثابت² را پیشنهاد کردند. به طور خاص، آنها از این واقعیت استفاده می کنند که هنگام در نظر گرفتن همه علل مستقیم یک متغیر هدف، توزیع شرطی هدف به شرط علل مستقیم، هنگام تداخل با همه متغیرهای دیگر در مدل به جز خود هدف، تغییر نخواهد کرد.

اگرچه *ICP* اولین تلاش برای ارتباط دادن روش ثابت³ به علیت است، اما محدودیت های متعددی دارد. ابتدایی ترین آنها، الزامات سختگیرانه برای ناهمگنی⁴ است زیرا قدرت *ICP* به شدت به کیفیت محیط های موجود (یا اغتشاشات⁵) بستگی دارد. اگر زیرجمعیت های آشفته موجود⁶ کافی نباشد یا حتی یک محیط واحد وجود نداشته باشد، کارایی *ICP* از بین خواهد رفت.

2-3-3- یادگیری ثابت

طبق روش های مبتنی بر استنتاج علی، روش های یادگیری ثابت، که با حداقل سازی ریسک ثابت⁷ مشخص می شوند، مکانیسم های علی پنهان⁸ را هدف قرار می دهند و *ICP* را به حالت کاربردی تر و عمومی تر گسترش می دهند. در حالی که روش های پیش بینی علی، سطح متغیر خام⁹ را در نظر می گیرند.

بر اساس *IRM* که با بازنمایی $\Phi(X)$ به گونه ای است که $E[Y | \Phi(X)]$ ثابت می ماند، کارهای بعدی تغییراتی را در این راستا به صورت قاعده مندتر از فرض عدم تغییر *IRM*¹⁰ پیشنهاد کرده اند که منجر به جایگزین های مشابه می شود. آهوجا و همکاران [7] از تئوری بازی ها در این مسئله استفاده کرده اند و طبقه بندی کننده خطی در *IRM* را با مجموعه ای از طبقه بندی کننده ها از محیط های مختلف جایگزین می کنند. جین و همکاران [7] تنظیم کننده *IRM*¹¹ را با تاثیر پیش بینی کننده¹ جایگزین می کنند و

¹ invariance

² ICP- Invariant Causal Prediction

³ invariance

⁴ heterogeneity

⁵ perturbations

⁶ available perturbed subpopulations

⁷ IRM- invariant risk minimization

⁸ latent causal mechanisms

⁹ raw variable level

¹⁰ IRM's invariance assumption

¹¹ regularizer of IRM

محدودیت های سختگیرانه تری را بر $\Phi(X)$ تحمیل می کنند. کروگ و همکاران [7] پیشنهاد می کنند که واریانس ریسک ها در بین محیط ها جریمه شوند، در حالی که Xie و همکاران. [7] تقریباً همان هدف را مطرح می کنند اما جریمه اصلی را با جذر واریانس جایگزین می کنند. ماهاجان و همکاران [7] یک قاعده ساز کنتراست² را معرفی می کند که بازنمایی اشیاء مشابه را در بین محیط ها مطابقت می دهد. و کریگر و همکاران مشکل از دست رفتن برچسب های محیطی IRM را هدف قرار می دهند و استنتاج محیطی برای یادگیری ثابت³ را برای یادگیری محیط های IRM که جریمه IRM را به حداکثر می رسانند، افزایش می دهند. کل الگوریتم دو مرحله ای است که ابتدا محیط ها را طبق یک مدل مرجع مغرضانه⁴ تولید می کند و سپس یادگیری ثابت را با محیط های یادگیری انجام می دهد.

در حالی که نتایج در IRM امیدوارکننده به نظر می رسد، روزنفلد و همکاران [7] به برخی از مشکلات IRM در مسائل طبقه بندی اشاره می کنند. در حالت خطی، آن ها شرایط ساده ای را ارائه می کنند که در آن راه حل بهینه اغلب نمی تواند پیش بینی کننده ثابت بهینه⁵ را بازیابی کند. به طور خاص، روزنفلد و همکاران [7] نشان می دهند که یک راه حل امکان پذیر وجود دارد که فقط از ویژگی های محیطی استفاده می کند و در عین حال عملکرد بهتری نسبت به پیش بینی کننده ثابت بهینه⁶ در همه محیط ها دارد. در یک حالت غیر خطی، آنها نشان می دهند که IRM می تواند به طور فاجعه باری شکست بخورد، مگر اینکه داده های آزمون به اندازه کافی شبیه توزیع آموزش باشند.

3-3-3- یادگیری پایدار⁷

در مقایسه با تعمیم دامن و یادگیری علی، یادگیری پایدار راه دیگری را برای ترکیب استنتاج علی با یادگیری ماشین ایجاد می کند، که به طور قابل توجهی الزامات محیط های متعدد را تسهیل می کند.

مسئله 2 (تنظیمات یادگیری پایدار): با توجه به داده های آموزشی $D^e = (X^e, Y^e)$ از یک محیط $e \in \text{supp}(\varepsilon_{all})$ ، هدف یادگیری پایدار، یادگیری یک مدل پیش بینی با عملکرد یکنواخت مناسب در

¹ predictive regret

² contrastive regularizer

³ EIIL- Environment Inference for Invariant Learning

⁴ a biased reference model

⁵ optimal invariant predictor

⁶ the optimal invariant predictor

⁷ Stable Learning

هر محیط ممکن در $supp(\varepsilon_{all})$ است.

برای حل چنین مشکل دشواری، با کمک استراتژی های متعادل کننده متغیر¹، شن و همکاران [12] پیشنهاد می کنند که همه متغیرها را به عنوان رفتار² در نظر بگیریم و مجموعه ای از وزن های نمونه سراسری را یاد بگیریم که می تواند سوگیری مخدوش کننده³ را برای همه رفتار های بالقوه⁴ از توزیع داده ها حذف کند. آنها یک تابع زیان متعادل سراسری⁵ را به دست می آورند که می تواند به راحتی به مسائل استاندارد یادگیری ماشین به عنوان تنظیم کننده⁶ متصل شود، همانطور که در معادله 3 نشان داده شده است:

$$\sum_{j=1}^p \left\| \frac{X_{-j}^T \cdot (W \odot I_j)}{W^T \cdot I_j} - \frac{X_{-j}^T \cdot (W \odot (1 - I_j))}{W^T \cdot (1 - I_j)} \right\|_2^2, \quad (3)$$

که W نشان دهنده ی وزن های نمونه و $\left\| \frac{X_{-j}^T \cdot (W \odot I_j)}{W^T \cdot I_j} - \frac{X_{-j}^T \cdot (W \odot (1 - I_j))}{W^T \cdot (1 - I_j)} \right\|_2^2$ نشان دهنده زیان تعادل مخدوش کننده⁷ هنگام تنظیم ویژگی j به عنوان متغیر رفتار⁸ است، و X_{-j} ویژگی های باقی مانده (یعنی مخدوش کننده ها) به جز ستون j ام است. I_j به معنای ستون j ام I است و I_{ij} به وضعیت رفتار واحد i هنگام تنظیم ویژگی j به عنوان متغیر رفتار اشاره دارد. با به حداقل رساندن زیان تعادل سراسری⁹، سوگیری مخدوش کننده¹⁰ را می توان در مقیاس سراسری حذف کرد.

علاوه بر این، کوانگ و همکاران. [7] بازنمایی ویژگی بدون نظارت را در مرحله تعادل سراسری¹¹ با رمزگذار های خودکار¹² [7] ترکیب می کنند و تنظیم کننده اصلی¹³ را به نسخه "عمیق" به عنوان معادله 4 تغییر

¹ variable balancing strategies

² treatment

³ confounding bias

⁴ potential treatments

⁵ global balancing loss

⁶ regularizer

⁷ loss of confounder balancing

⁸ treatment variable

⁹ global balancing loss

¹⁰ the confounding bias

¹¹ the global balancing stage

¹² auto-encoders

¹³ original regularizer

می دهند:

$$\sum_{j=1}^p \left\| \frac{\phi(X_{:, -j})^T \cdot (W \odot X_{:, j})}{W^T \cdot X_{:, j}} - \frac{\phi(X_{:, -j})^T \cdot (W \odot (1 - X_{:, j}))}{W^T \cdot (1 - X_{:, j})} \right\|_2^2, \quad (4)$$

روش های فوق به ویژگی های باینری محدود می شوند، زیرا مسائل مطرح شده در مورد استنتاج علی در محدوده رفتار باینری هستند. هنگامی که متغیر رفتار¹ قطعی یا پیوسته است، روش های متعادل سازی قدیمی² دیگر امکان پذیر نیستند، زیرا سطح رفتار³ می تواند بی نهایت بزرگ باشد. برای کاهش چنین محدودیت هایی و پرداختن به مشکل رفتار پیوسته⁴، کوانگ و همکاران. [7] پیشنهاد به یادگیری مجموعه ای از وزن های نمونه می کنند تا توزیع وزنی رفتار و مخدوش کننده⁵ بتواند شرایط استقلال را برآورده کند، که با این واقعیت مطابقت دارد که اگر رفتار و مخدوش کننده⁷ مستقل باشند، اثر رفتار را می توان به طور دقیق تخمین زد.

4-3- بهینه سازی برای تعمیم خارج از توزیع

برای پرداختن به مسئله تعمیم خارج از توزیع، جدا از یادگیری بازنمایی بدون نظارت و مدل های یادگیری انتها به انتها، روش های بهینه سازی با ضمانت های نظری، اخیراً به صورت گسترده مورد توجه قرار گرفته اند که هم به صورت آگنوستیک مدل⁶ و هم آگنوستیک ساختار داده⁷ می باشند. در این بخش ابتدا هدف این روش های بهینه سازی خارج از توزیع را معرفی می کنیم و سپس روش های بهینه سازی از جمله بهینه سازی مقاوم توزیعی⁸ و بهینه سازی مبتنی بر ثابت⁹ را بررسی می کنیم.

به منظور پرداختن به مسئله تعمیم خارج از توزیع از منظر بهینه سازی، مسئله تعمیم خارج از توزیع برای کنترل خطای پیش بینی بدترین حالت در بین \mathcal{E}_{all} فرموله شده است که به شکل زیر می باشد:

-
- ¹ treatment variable
 - ² traditional balancing methods
 - ³ treatment level
 - ⁴ the continuous treatment problem
 - ⁵ treatment and confounder
 - ⁶ model agnostic
 - ⁷ data structure agnostic
 - ⁸ Distributionally Robust Optimization
 - ⁹ Invariant-Based Optimization

$$\operatorname{argmin}_f \max_{e \in \operatorname{supp}(\varepsilon_{all})} \mathcal{L}(f|e) \quad (5)$$

که در آن ε_{all} متغیر تصادفی در همه محیط های ممکن است، و برای همه $e \in \operatorname{supp}(\varepsilon_{all})$ ، توزیع داده ها و برچسب $P^e(X, Y)$ می تواند کاملاً متفاوت از توزیع آموزشی $P_{tr}(X, Y)$ باشد. $\mathcal{L}(f|e) = \mathbb{E}[l(f(X), Y)|e] = \mathbb{E}^e[l(f(X^e), Y^e)]$ ریسک پیش بینی کننده f در محیط e و $l(.,.): y \times y \rightarrow \mathbb{R}^+$ تابع زیان است. به طور شهودی، هدف روش های بهینه سازی تضمین عملکرد بدترین حالت تحت تغییرات توزیعی است.

3-4-1- بهینه سازی مقاوم توزیعی

بهینه سازی مقاوم توزیعی¹، مستقیماً مسئله تعمیم خارج از توزیع را با بهینه سازی خطای بدترین حالت در یک مجموعه توزیع عدم قطعیت² حل می کند تا از مدل در برابر تغییرات توزیعی بالقوه در مجموعه عدم قطعیت محافظت شود. اغلب توسط شرایط لحظه ای یا پشتیبانی³، واگرایی- f ⁴ و فاصله $Wasserstein$ ، محدود می شود. تابع هدف روش های DRO را می توان به صورت زیر خلاصه کرد:

$$\operatorname{argmin}_f \sup_{Q \in \mathcal{P}(P_{tr})} \mathbb{E}_{X,Y \sim Q}[l(f(X), Y)] \quad (6)$$

که در آن $\mathcal{P}(P_{tr})$ مجموعه توزیعی است که در پیرامون توزیع آموزشی P_{tr} قرار دارد و $l(.,.): y \times y \rightarrow \mathbb{R}^+$ تابع زیان است. روش های مختلف DRO انواع مختلفی از محدودیت ها را برای فرمول بندی مجموعه توزیع $\mathcal{P}(P_{tr})$ و الگوریتم بهینه سازی مختلف اتخاذ می کنند.

3-4-2- بهینه سازی مبتنی بر عدم تغییر⁵

برخلاف روش های DRO که مستقیماً برای بدترین حالت بدون هیچ فرضی اضافی بهینه سازی را انجام می دهند، روش های بهینه سازی مبتنی بر عدم تغییر خاصیت عدم تغییر را در داده ها فرض می کنند و از

¹ DRO-Distributionally Robust Optimization

² uncertainty distribution set

³ moment or support conditions

⁴ f-divergence

⁵ Invariance-Based Optimization

محیط های متعدد برای یافتن چنین تغییر ناپذیری برای تعمیم تحت تغییرات توزیعی استفاده می کنند. چانگ و همکاران [7] و کویاما و همکاران [7] بازنمایی عدم تغییر مورد نظر را با استفاده از تئوری اطلاعات، فرموله کرده اند و پیشنهاد به پیدا کردن پیش بینی کننده تغییر ناپذیر حداکثری¹ در محیط های یادگیری برای کنترل بدترین خطا در معادله 5 می کنند.

3-5- مجموعه داده ها و معیارهای ارزیابی

برای توسعه و کمک به تحقیقاتی که در راستای تعمیم خارج از توزیع انجام می شوند، ارزیابی عملکرد تعمیم خارج از توزیع الگوریتم های مختلف به صورت منطقی و دقیق، از اهمیت زیادی برخوردار است. ارزیابی یک الگوریتم معمولاً از دو بخش مجزا به اسم مجموعه داده ها و معیارهای ارزیابی تشکیل شده است. در این بخش، مجموعه داده ها و معیارهای ارزیابی را که معمولاً به عنوان معیار خارج از توزیع² در تحقیقات و مقالات موجود استفاده می شوند را ارائه می کنیم.

3-5-1- مجموعه داده ها³

مجموعه داده ها را می توان بر اساس معیارهای مختلفی طبقه بندی کرد، برای مثال، داده های ساختگی و داده های دنیای واقعی، داده های ویژگی خام و داده های پیچیده. زمینه های تحقیقاتی مختلف از انواع مختلف مجموعه داده ها استفاده می کنند، برای مثال، یادگیری آماری ای که به صورت مرسوم استفاده می شود اغلب از داده های ساختگی⁴ استفاده می کند. بینایی کامپیوتر اغلب از داده های تصویری در دنیای واقعی استفاده می کند. در مورد مسئله تعمیم خارج از توزیع، متفاوت از وظایف یادگیری ماشین سنتی که مبتنی بر فرض *i.i.d* هستند، لازم است که تغییرات توزیعی را برای شبیه سازی توزیع داده های آزمون که ناشناخته هستند به کار بگیریم یا ایجاد کنیم، تا آزمایش شود که آیا یک الگوریتم می تواند به توزیع های دیده نشده تعمیم یابد یا خیر. بنابراین، مطابق با کارهای اخیر در مورد تعمیم خارج از توزیع، استفاده از داده های مصنوعی و ساده و داده های دنیای واقعی و پیچیده برای تأیید اثربخشی روش تعمیم خارج از توزیع ضروری می باشد.

¹ MIP- Maximal Invariant Predictor

² OOD benchmarks

³ Datasets

⁴ synthetic data

1-1-5-3- داده های ساختگی

داده های ساختگی¹ ابزار مهمی برای شبیه سازی تغییرات توزیعی قابل توضیح و کنترل هستند. همانطور که اوبین و همکاران [13] دریافتند که روش های تعمیم خارج از توزیع اخیر در برخی از مسائل خطی ساده کم بُعد ضعیف عمل می کنند، لازم است روش های تعمیم خارج از توزیع را روی چنین داده های ساده اما چالش برانگیزی آزمایش کنیم، که می تواند نشان دهد آیا و تا چه حد یک الگوریتم می تواند در برابر نوع خاصی از تغییرات توزیعی مقاوم باشد.²

به طور کلی، سه مکانیسم برای شبیه سازی تغییرات توزیعی در محیط ها وجود دارد، به نام های سوگیری انتخاب³، سوگیری مخدوش کننده⁴ و اثر ضد علی⁵، که با آن ها می توان انواع خاصی از تغییرات توزیعی را به درجات مختلف شبیه سازی کرد و به وضوح اثر واقعی الگوریتم را توجیه کرد. در این مکانیسم ها، *covariate* های X به دو گروه $X = [S, V]^T$ ، مربوط به بخش های ثابت و متغیر داخل داده ها تقسیم می شوند. و فرض بر این است که $P(Y|S)$ در سراسر محیط ها ثابت می ماند، و $P(Y|V)$ با مکانیسم های مختلف دچار تغییر می شود⁶، که تغییرات توزیعی را به همراه دارد.

الف) سوگیری انتخاب⁷

کوانگ و همکاران [7] یک مکانیسم سوگیری انتخاب را برای معرفی تغییرات توزیعی پیشنهاد می کنند. در الگوریتم آن ها، $P(Y|V)$ توسط سوگیری انتخاب دچار تغییر می شود.

ب) سوگیری مخدوش کننده⁸

سوگیری مخدوش کننده نیز یکی از رایج ترین ابزارهای تغییرات توزیعی است. در مقایسه با سوگیری انتخاب، در این روش، ویژگی متغیر V به هدف Y با در نظر گرفتن مخدوش کننده مشاهده نشده C^9 مربوط می شود.

¹ Synthetic data

² resist

³ selection bias

⁴ confounding bias

⁵ anti-causal effect

⁶ perturbed

⁷ selection bias

⁸ confounding bias

⁹ unobserved confounder C

ج) اثر ضد علی¹

علاوه بر مکانیسم های فوق، آریوفسکی و همکاران [7] و لیو و همکاران [7] یک مکانیسم ضد علت را برای تغییر $P(Y|V)$ معرفی می کنند.

2-1-5-3- داده های واقعی²

داده های مصنوعی به راحتی می توانند تغییرات توزیعی را با درجات مختلف ایجاد کنند، همچنین می توانند یک الگوریتم را به صورت کامل بررسی کنند که آیا می تواند به توزیع های دیده نشده تعمیم یابد یا خیر. با این حال، داده های مصنوعی نسبتاً ساده هستند و تولید داده های پیچیده (مثلاً تصاویر) دشوار است. علاوه بر این، اینکه آیا یک الگوریتم می تواند مسائل تغییر توزیعی در دنیای واقعی را حل کند نیز یک معیار مهم برای ارزیابی یک روش تعمیم خارج از توزیع است. بنابراین، لازم است که مجموعه داده های دنیای واقعی را در نظر بگیریم. در اینجا، ما معروف ترین مجموعه داده های دنیای واقعی مورد استفاده در تحقیقات خارج از توزیع را بررسی می کنیم. از جمله مجموعه داده های تصویری و سایر فرم های داده (به عنوان مثال، داده های جدولی، داده های زبان). خلاصه ای از این مجموعه داده ها در جدول 1 آمده است.

Image Data Set	ImageNet-Variant [165], [166], [167]	Colored MNIST [2]	MNIST-R [169]	Waterbirds [153]	Camelyon17 [173]	VLCS [172]	PACS [170]
# Domains	-	3	6	2	5	4	4
# Categories	-	2	10	2	2	5	7
# Examples	-	-	6,000	4,800	45w	2,800	10w
Shift Type	Adversarial Policy	Color	Angle	Background	Hospital	Data Source	Style
Image Type	Mixed Type	Digits	Digits	Birds	Tissue Slides	Real Objects	Mixed Type
Image Data Set	Office-Home [171]	DomainNet [174]	iWildCam [175]	FMoW [178]	PovertyMap [179]	NICO [177]	
# Domains	4	6	323	16 × 5	23 × 2	188	
# Categories	65	345	182	62	Real Value	19	
# Examples	15w	50w	20w	50w	2w	2.5w	
Shift Type	Style	Style	Location	Time, Location	Country, Urban/Rural	Background, Attribute, Action, View and Co-occurring Object	
Image Type	Mixed Type	Mixed Type	Real Animals	Satellite	Satellite	Real Objects	

جدول 1: مجموعه داده های تصویری که معمولاً برای تعمیم OOD استفاده می شوند. نوع تغییر ($Shift type$)

نشان دهنده نوع تغییرات توزیعی است و نوع مختلط ($mixed type$) در نوع تصویر ($image type$) به این

معنی است که هم تصاویر واقعی و هم غیر واقعی وجود دارد.

مجموعه داده های تصویری

با توسعه سریع بینایی کامپیوتر، تعدادی از مجموعه داده های تصویری ایجاد شده است. ما آن ها را با توجه

¹ anti-causal effect

² Real-World Data

به میزان شبیه سازی تغییرات توزیعی به سه دسته طبقه بندی می کنیم:

1. داده های تبدیل مصنوعی¹
2. داده های رام نشده ثابت²
3. داده های رام نشده قابل کنترل³

الف) داده های تبدیل مصنوعی¹

اگرچه بیشتر مجموعه داده های تصویر برای تعمیم خارج از توزیع تولید نمی شوند، اما محققان آن ها را با برخی تبدیل های مصنوعی تغییر می دهند تا تغییرات توزیعی را شبیه سازی کنند. رایج ترین آن ها، شامل انواع *ImageNet* (به عنوان مثال *ImageNet-A*, *ImageNet-C*, *ImageNet-R*) است که خط مشی انتخاب داده یا تغییرات⁴ خاصی را برای تولید داده های آزمایشی با تغییرات توزیعی اتخاذ می کنند. بقیه مجموعه داده ها روی *MNIST* کار کرده اند که انواع *MNIST* را به وجود آورده اند (مانند *MNIST* رنگی، *MNIST-R*)، که محیط های مختلف را با رنگ آمیزی یا چرخش تصاویر اصلی شبیه سازی می کنند. از آنجایی که این مجموعه داده ها به خوبی طراحی شده اند، آنها را برای مطالعه اولیه و تأیید اثربخشی الگوریتم های تعمیم خارج از توزیع استفاده می کنند.

ب) داده های رام نشده ثابت⁵

چند مجموعه داده برای اعتبار سنجی⁶ تعمیم خارج از توزیع ساخته شده اند که عمدتاً از پس زمینه ها یا محیط های دنیای واقعی استفاده می کنند. *PACS* و *Office-Home* که به طور گسترده در تعمیم دامنه استفاده می شود، از سبک تصویر (به عنوان مثال، هنری، کارتونی) برای متمایز کردن دامنه ها/توزیع ها استفاده می کنند، و *VLCS* داده ها را از چهار منبع مستقل جمع آوری می کند. علاوه بر این، *Camelyon17* حاوی اسلایدهای بافت نمونه برداری شده است که روی آنها پس-پردازش⁷ انجام شده و از بیمارستان های مختلف جمع آوری شده است. *PACS* *DomainNet* را در مقیاسی بسیار بزرگتر، متشکل از حوزه ها و دسته های بیشتری گسترش می دهد.

¹ synthetic transformation data

² fixed wild data

³ controllable wild data

⁴ perturbations

⁵ fixed wild data

⁶ validation

⁷ post-processed

ج) داده های رام نشده قابل کنترل¹

اخیراً، مجموعه داده ای وجود دارد که روش های انعطاف پذیرتر و قابل کنترل تری را برای شبیه سازی تغییرات توزیعی امکان پذیر می سازد که با *NICO* مشخص می شود. روش *NICO* به طور قابل قبولی زمینه های بصری² را با تنظیمات مختلف از جمله پس زمینه، ویژگی³، نما⁴ و غیره انتخاب می کند و محیط های مختلفی را تولید می کند. این روش می تواند انواع مختلفی از تغییرات توزیعی واقعی با زمینه های متنوع را شبیه سازی کند، و با اندازه نمونه یکسان در هر زمینه، می توان درجات مختلفی از تغییرات توزیعی را به راحتی به دست آورد. این دو ویژگی باعث می شود روش *NICO* تنظیمات انعطاف پذیر تغییرات توزیعی را دارا باشد.

2-5-3- معیار های ارزیابی

علاوه بر مجموعه داده ها، معیارهای ارزیابی نیز برای ارزیابی مناسب الگوریتم های تعمیم خارج از توزیع مهم هستند که باید در نظر گرفته شوند. در مقایسه با مسئله ی *iid* که در آن تنها یک توزیع آزمون در نظر گرفته می شود، اغلب، توزیع های آزمون متعددی در مسائل تعمیم خارج از توزیع وجود دارند تا توزیع های دیده نشده را بهتر به تصویر بکشند. علاوه بر این، *Ye* و همکاران [14] به طور تجربی دریافتند که دقت آزمون روی یک محیط واحد، ارزیابی الگوریتم ها را در یک سناریوی خارج از توزیع به درستی انجام نمی دهد. بنابراین، برای ارزیابی مناسب الگوریتم های تعمیم خارج از توزیع، اطلاعات آماری بیشتری را باید در مورد دقت محیط های آزمون مختلف در نظر گرفت. در اینجا، ما سه معیار ارزیابی، از جمله دقت متوسط، دقت در بدترین حالت، و انحراف معیار دقت را بررسی می کنیم. برای راحتی، دقت مدل در K محیط آزمون را به ترتیب $\{acc_1, \dots, acc_k\}$ فرض می کنیم.

الف) دقت متوسط

دقت متوسط \overline{Acc} روی توزیع های آزمون، ساده ترین راه برای ارزیابی اثربخشی الگوریتم های خارج از توزیع است که معمولاً در مقاله های تعمیم خارج از توزیع استفاده می شود و به صورت زیر محاسبه می شود:

¹ controllable wild data

² visual contexts

³ attribute

⁴ view

$$\overline{Acc} = \frac{1}{k} \sum_{k=1}^k acc_k, \quad (7)$$

دقت متوسط، عملکرد کلی را در میان توزیع های آزمون اندازه گیری می کند، اما نمی تواند نوسانات عملکرد یک الگوریتم را توصیف کند. علاوه بر این، میانگین دقت به طور یکسان با همه توزیع های آزمون بدون در نظر گرفتن ویژگی هر یک، رفتار می کند که ممکن است از توزیع هایی که مکررا رخ می دهند فریب بخورد.

ب) دقت در بدترین حالت

این دقت به طور گسترده در مقالات *DRO* استفاده می شود. Acc_{worst} به عنوان دقت در بدترین حالت روی توزیع های آزمون تعریف می شود:

$$Acc_{worst} = \min_{k \in [k]} acc_k \quad (8)$$

دقت در بدترین حالت، قابلیت اطمینان یک الگوریتم را نشان می دهد.

ج) انحراف معیار¹

انحراف معیار دقت روی توزیع های آزمون (Acc_{std})، تغییرات عملکرد در توزیع های مختلف را اندازه گیری می کند که به صورت زیر تعریف می شود:

$$Acc_{std} = \sqrt{\frac{1}{k-1} \sum_{k=1}^k (acc_k - \overline{Acc})^2} \quad (9)$$

این متریک حساسیت یک الگوریتم را اندازه گیری می کند و استحکام² و پایداری³ الگوریتم را نشان می دهد

¹ STD; Standard Deviation

² robustness

³ stability

که این موضوع برای یک الگوریتم تعمیم خارج از توزیع مهم است.

3-6- نتیجه‌گیری

مسئله تعمیم خارج از توزیع اخیراً بسیار مورد توجه قرار گرفته است و همانطور که می‌دانیم برای استقرار الگوریتم‌های یادگیری ماشین حیاتی است. در این فصل، تعریف، شاخه‌های اصلی روش‌ها، و مجموعه داده‌ها و معیارهای ارزیابی مسئله تعمیم خارج از توزیع را بررسی کردیم. بر اساس تجزیه و تحلیل، ما با چندین چالش بالقوه روبرو شدیم که در فصل بعد آن‌ها را مطرح خواهیم کرد.

فصل 4:

نتیجه‌گیری و کارهای آینده

1-4- نیاز به محیط‌های متعدد

اکثر روش‌های تعمیم خارج از توزیع به چندین محیط آموزشی نیاز دارند. با این حال، مجموعه داده‌های مدرن اغلب با ادغام داده‌ها از منابع متعدد بدون حفظ برچسب‌های منبع جمع‌آوری می‌شوند، که به شدت استقرار روش‌های تعمیم خارج از توزیع را در سناریوهای واقعی محدود می‌کند. بنابراین، عملی‌تر و واقع‌بینانه‌تر است که ما فقط به یک محیط آموزشی با ناهمگونی نهفته دسترسی داشته باشیم. اخیراً، اگرچه برخی از مقالات [15] وجود دارند که سعی می‌کنند از ناهمگونی پنهان استفاده کنند و نیاز برای محیط‌های متعدد را کاهش دهند، اما نحوه جست‌وجو و استفاده مناسب از ناهمگونی پنهان درون داده‌ها برای استقرار روش‌های تعمیم خارج از توزیع حیاتی است. و باید به عنوان ادامه‌ی کارهای آینده در نظر گرفته شود.

2-4- ارزیابی‌های منطقی

اگرچه معیارهای ارزیابی از جمله داده‌های آزمون و مکانیسم‌های انتخاب مدل، برای الگوریتم‌های یادگیری ماشین کلاسیک تحت فرض *iid* به خوبی توسعه داده شده‌اند اما آنها را نمی‌توان مستقیماً در سناریوهای خارج از توزیع استفاده کرد. از آنجایی که توزیع آزمون ناشناخته و متفاوت از آموزش است، نحوه طراحی تنظیمات تجربی منطقی و واقعی همچنان یک مشکل چالش‌برانگیز است. علاوه بر این، مکانیسم انتخاب مدل نیز مهم است، زیرا انتخاب داده‌های اعتبارسنجی در سناریوهای خارج از توزیع مهم می‌باشد. تحقیقات نشان می‌دهند که الگوریتم‌های تعمیم دامنه بدون استراتژی انتخاب مدل ناقص هستند. همچنین بیان می‌کنند که اثرات واقعی بسیاری از روش‌های تعمیم دامنه ضعیف است که نشان می‌دهد معیارهای ارزیابی موجود برای اعتبارسنجی الگوریتم‌های تعمیم خارج از توزیع ناکافی هستند. بنابراین، ضروری است که معیارهای ارزیابی معقول‌تری برای تعمیم خارج از توزیع ایجاد شوند.

3-4- معرفی موضوع مورد نظر برای پایان‌نامه

به نظر بنده مدل‌های موجود هنوز به خوبی بازنمایی‌های مناسبی برای تعمیم خارج از توزیع بدست نمی‌آوردند. بنده تصمیم دارم همچون مقاله‌ی CIRL که در فصل دوم به آن پرداخته شد یک SCM مناسب‌تر برای یادگیری بازنمایی‌های ثابت از طریق یادگیری علی ارائه دهم که بتواند تعمیم مناسبی روی دامنه هدف دیده نشده داشته باشد.

مراجع

- [1] J. Peters, D. Janzing, and B Schölkopf. Elements of causal inference - foundations and learning algorithms. MIT Press, Cambridge, MA, USA, 2017.
- [2] H. Reichenbach. The Direction of Time. University of California Press, Berkeley, CA,, 1956.
- [3] Bernhard Schölkopf, Dominik Janzing, Jonas Peters, Eleni Sgouritsa, Kun Zhang, and Joris M. Mooij. On causal and anticausal learning. In ICML, 2012.
- [4] F. Lv, J. Liang, S. Li, B. Zang. Causality Inspired Representation Learning for Domain Generalization. CVPR, 2022.
- [5] F. Locatello, S. Bauer, M. Lucic, G. Raetsch, S. Gelly, B. Scholkopf, " and O. Bachem, "Challenging common assumptions in the unsupervised learning of disentangled representations," in international conference on machine learning. PMLR, 2019, pp. 4114–4124.
- [6] I. Albuquerque, J. Monteiro, M. Darvishi, T. H. Falk, and I. Mitliagkas, "Adversarial target-invariant representation learning for domain generalization," 2020.
- [7] Z. Shen, J. Liu, Y. He, X. Zhang. Towards Out-Of-Distribution Generalization: A Survey. Arxiv, 2021.
- [8] X. Yue, Y. Zhang, S. Zhao, A. Sangiovanni-Vincentelli, K. Keutzer, and B. Gong, "Domain randomization and pyramid consistency: Simulation-to-real generalization without accessing target domain data," in Proceedings of the IEEE/CVF International Conference on Computer Vision, 2019, pp. 2100–2110.
- [9] K. Zhou, Y. Yang, T. Hospedales, and T. Xiang, "Deep domainadversarial image generation for domain generalisation," in Proceedings of the AAAI Conference on Artificial Intelligence, vol. 34, no. 07, 2020, pp. 13 025–13 032.
- [10] F. Qiao, L. Zhao, and X. Peng, "Learning to learn single domain generalization," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 12 556–12 565.
- [11] B. Schölkopf, F. Locatello, S. Bauer. Towards Causal Representation Learning. arxiv, 2021.
- [12] Z. Shen, P. Cui, K. Kuang, B. Li, and P. Chen, "Causally regularized learning with agnostic data selection bias." in ACM Multimedia, 2018, pp. 411–419.
- [13] B. Aubin, A. Słowiak, M. Arjovsky, L. Bottou, and D. LopezPaz, "Linear unit-tests for invariance discovery," arXiv preprint arXiv:2102.10867, 2021.
- [14] H. Ye, C. Xie, Y. Liu, and Z. Li, "Out-of-distribution generalization analysis via influence function," arXiv preprint arXiv:2101.08521, 2021.
- [15] E. Creager, J.-H. Jacobsen, and R. Zemel, "Environment inference for invariant learning," in International Conference on Machine Learning. PMLR, 2021, pp. 2189–2200.

واژه نامه

واژه نامه فارسی به انگلیسی

Interventions	تغییرات
Confounders	عوامل مخدوش کننده
Stable Learning	یادگیری پایدار
Selection bias	سوگیری انتخاب
Meta-Learning	فرا-یادگیری
Invariant Representation Learning	یادگیری بازنمایی ثابت
Domain Augmentation	تقویت دامنه
parametric	دارای پارامتر
ground-truth	برچسب واقعی
stream	جریان
shift	تغییر
covariate shift	تغییر متغیر
SCM; structural causal model	مدل علّی ساختاری
category	دسته
causal factors	عوامل علّی
non-causal factors	عوامل غیر علّی
invariant causal mechanisms	مکانیسم های علّی ثابت
causal intervention	تغییرات علّی
category-related information	اطلاعات مربوط به هر دسته
CIRL; Causality Inspired Representation Learning	یادگیری بازنمایی الهام گرفته از علیت
perturbed domain-related information	با اطلاعات مرتبط با دامنه‌ی دچار تغییر شده
adversarial mask module	ماژول ماسک متخاصم
masker	پوشش‌دهنده
domain-invariant marginally or conditionally	به صورت حاشیه ای یا شرطی ثابت-دامنه
Common Cause Principle	اصل علت مشترک
cross-entropy loss	تابع زیان آنتروپی متقابل

<i>ICM; Independent Causal Mechanisms</i>	اصل مکانیزم های علی مستقل
<i>causal factorization</i>	فاکتورسازی علی
<i>batch size</i>	اندازه دسته ها
<i>inferior dimensions</i>	ابعاد با ارزش پایین تر
<i>superior dimensions</i>	ابعاد با ارزش بالاتر
<i>masked</i>	ماسک شده
<i>inference</i>	استنتاج
<i>inductive bias</i>	سوگیری استقرایی
<i>Disentangled Representation Learning</i>	یادگیری بازنمایی تفکیک شده
<i>embed</i>	جاسازی
<i>Distributionally Robust Optimization</i>	بهینه سازی مقاوم توزیعی
<i>Invariance-Based Optimization</i>	بهینه سازی مبتنی بر عدم تغییر
<i>sparsity</i>	پراکندگی
<i>latent factor</i>	عوامل پنهان
<i>semi-supervised</i>	نیمه نظارتی
<i>end-to-end</i>	انتها به انتها
<i>DG; Domain Generalization</i>	تعمیم دامنه
<i>disturbance</i>	اختلال
<i>transferable</i>	قابل انتقال
<i>robust</i>	مقاوم
<i>domain adversarial learning</i>	یادگیری متخاصم دامنه
<i>domain alignment</i>	هم ترازی دامنه
<i>DANN</i>	شبکه عصبی متخاصم دامنه
<i>domain adaptation</i>	تطبیق دامنه
<i>manifold space</i>	فضای چندگانه
<i>class-specific adversarial networks</i>	شبکه های متخاصم کلاس خاص
<i>CIAN- conditional invariant adversarial network</i>	شبکه متخاصم ثابت شرطی
<i>MMD-maximum mean discrepancy</i>	حداکثر میانگین اختلاف

DICA- Domain-Invariant Component Analysis	تجزیه و تحلیل مؤلفه های ثابت دامنه
attribute regularization	تنظیم صفت
ensemble learning	یادگیری گروهی
MAML; model-agnostic meta-learning	فرایادگیری مدل-آگنوستیک
domain-specific	خاص-دامنه
DIUL; Domain-Irrelevant Unsupervised Learning	یادگیری بدون نظارت نامربوط دامنه
self-challenging dropout algorithm	الگوریتم حذف تصادفی خود چالش برانگیز
randomization based augmentation	تقویت مبتنی بر تصادفی بودن
ICP- Invariant Causal Prediction	پیش بینی علی ثابت
heterogeneity	ناهمگنی
available perturbed subpopulations	زیرجمعیت های آشفته موجود
raw variable level	سطح متغیر خام
predictive regret	تأثر پیش بینی کننده
contrastive regularizer	قاعده ساز کنتراست
EIIL- Environment Inference for Invariant Learning....	استنتاج محیطی برای یادگیری ثابت
biased reference model	مدل مرجع مغرضانه
variable balancing strategies	استراتژی های متعادل کننده متغیر
treatment	رفتار
potential treatments	رفتار های بالقوه
global balancing loss	زیان تعادل سراسری
treatment variable	متغیر رفتار
the global balancing stage	مرحله تعادل سراسری
auto-encoders	رمزگذار های خودکار
the continuous treatment problem	مشکل رفتار پیوسته
uncertainty distribution set	مجموعه توزیع عدم قطعیت
moment or support conditions	شرایط لحظه ای یا پشتیبانی
MIP- Maximal Invariant Predictor	پیش بینی کننده تغییر ناپذیر حداکثری
OOD benchmarks	معیار خارج از توزیع

synthetic data.....	داده های ساختگی.....
synthetic transformation data	داده های تبدیل مصنوعی.....
fixed wild data.....	داده های رام نشده ثابت.....
controllable wild data.....	داده های رام نشده قابل کنترل.....
visual contexts	زمینه های بصری.....

Abstract

Machine learning algorithms are usually built on the *iid* assumption that the training and test data are independent and identically distributed. This assumption means that the distribution of training data and test data are the same. In the real world, due to distributional shifts, this assumption is hardly fulfilled, which greatly reduces the accuracy of these classical machine learning algorithms. On the other hand, machine learning algorithms often use statistical models to model the dependence between data and labels, which aims to learn domain-independent representations. However, statistical models are superficial descriptions of reality because they only need to model dependencies rather than intrinsic causal mechanisms. When dependence changes with the target distribution, statistical models may fail to generalize. Causality, by focusing on the representation of structured knowledge about the data generation process that allows for interventions, can help to understand and address some of the limitations of current machine learning methods. Despite the success of statistical learning, these models provide a relatively superficial description of reality that holds only when the experimental conditions are invariant. Instead, the field of causal learning seeks to model the effect of distributional shifts with a combination of data-driven learning and assumptions not previously included in the statistical description of a system.

For the problem caused by the distribution shift, where the distribution of the test data is different from the training data, the out-of-distribution generalization problem arises through which the algorithm can generalize well on the unseen test data. In this seminar, we will review the methods of the out-of-distribution generalization problem. We will also explain one of its methods called causality, which has recently received extensive attention, in more detail.

Keywords: Representation Learning, Domain Generalization, Causal Inference, Out-Of-Distribution Generalization, Invariant Learning



IU ST

**Iran University of Science and Technology
School of Computer Engineering**

Causal Representation Learning for Out-Of-Distribution Generalization

**A Seminar Submitted in Partial Fulfillment of the Requirement for the
Degree of Master of Science in Computer Engineering – Artificial
Intelligence**

**By:
Hossein Rezaei**

**Supervisor:
Dr. Adel Rahmani**

Fall 2022