

Distribution Shifts And How We Can Learn from Them

Sean Kulinski

Probabilistic and Understandable Machine Learning Lab
- Purdue University -

View of the Next Hour:

- Part 1: Background on distribution shifts
 - What is distribution shift?
 - Why are distribution shifts (currently) problematic?
- Part 2: Detecting and understanding distribution shifts
 - So what are we doing about distribution shifts?
 - Looking forward, can we utilize distribution shifts to help us learn better?

About me :)

- Ph.D. Student in Computer Engineering @ Purdue University
- Belong to the Probabilistic and Understandable Machine Learning Lab lead by [Dr. David Inouye](#)
- Outside of research, I enjoy:
 - 🥾 hiking/backpacking,
 - 🚲 mountain biking
 - 🧠 spending time with friends and family
 - 🛠️ figuring out how things work



My partner and our dog



PUML Research Lab Group

My research path:

- My main research interest is:

“How can we build generalizable Machine Learning models for deployment to dynamic environments seen in the wild?”

Research Assistantship @ Purdue
(Lead By Dr. David Inouye)



Fall '19



Summer '19

Research Intern @ Lawrence Livermore National Lab
(Lead by Bhavya Kailkhura)

ML Scientist @ AbbVie
(Lead by David Masica)

abbvie

Fall '21 – Spring '22



Summer '22

Data Scientist Intern @ Microsoft365 Research
(Lead By Yi-Cheng Pan)

Research Intern @ Microsoft365 Research
(Lead By Ankur Mallick and Kevin Hu)



Summer '22

Part 1:

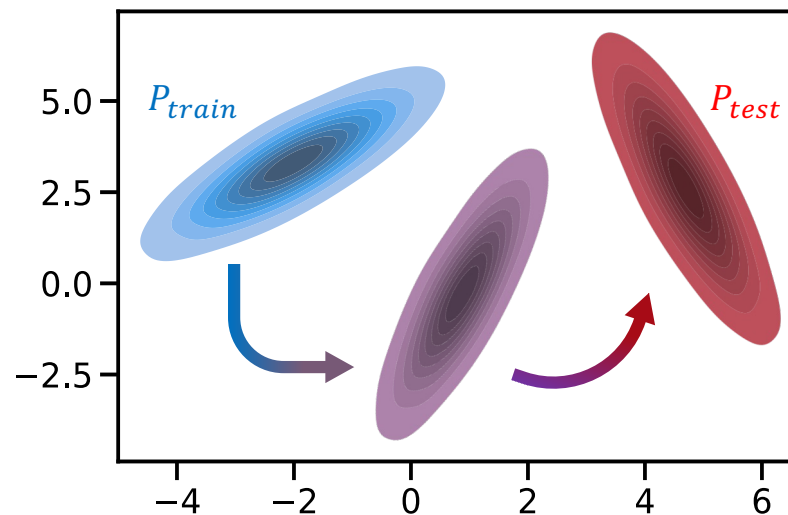
What are Distribution Shifts?

How the real world breaks fundamental ML assumptions.

A distribution shift is when a data distribution changes from what is expected

- In machine learning, a distribution shift is when a **testing distribution** no longer matches the **training distribution**

$$P_{test}(x) \neq P_{train}(x)$$



Most ML assumes train/test data distributions match

- Fundamental to most ML is the

i.i.d. assumption:

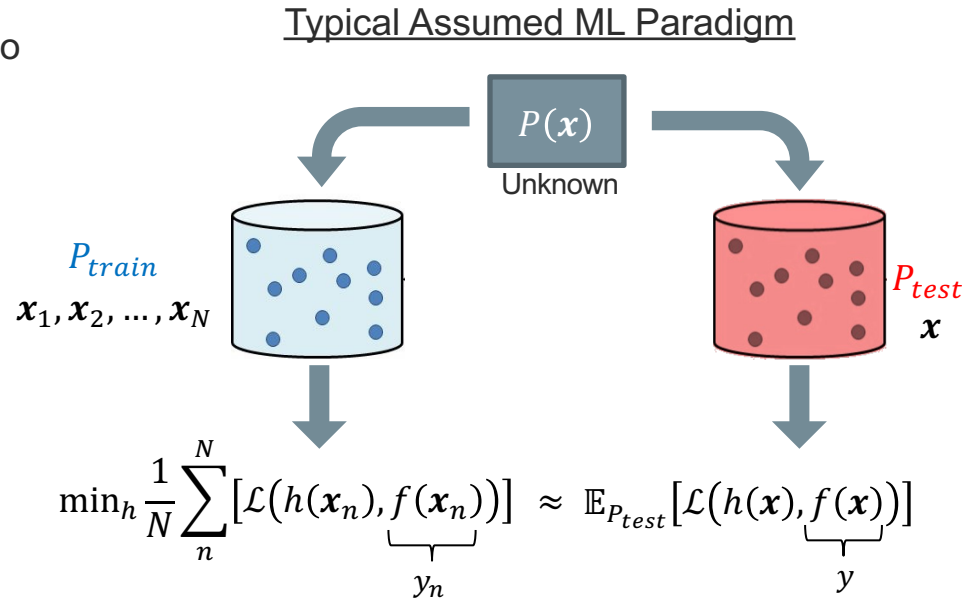
- Independent: All samples x are unrelated to each other

$$P(x_i | x_{i'}) = P(x_i) \quad \forall i \neq i'$$

- Identically Distributed: All samples x come from the same distribution

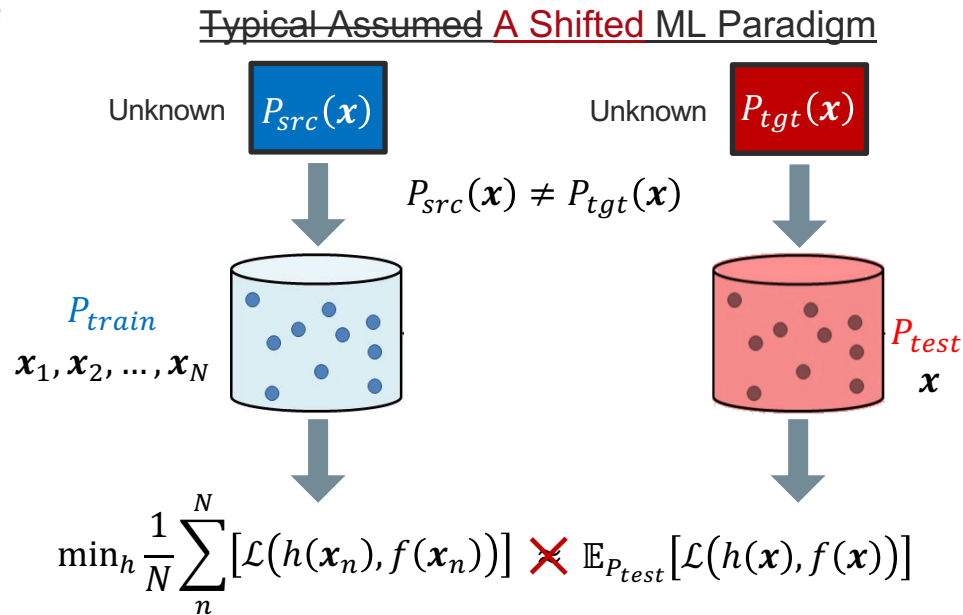
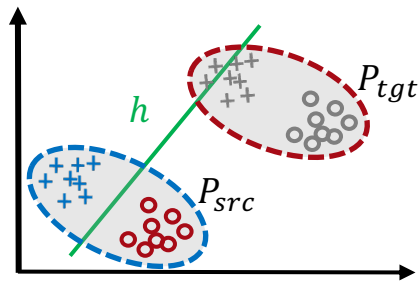
$$P_{train}(x) = P_{test}(x)$$

- The i.i.d. assumption allows our ML model h to generalize to P_{test}



Distribution shift violates this core assumption in ML

- Distribution shift usually breaks the identically distributed assumption
- Under distribution shift, the patterns learned by h might not hold under $P_{tgt}(\mathbf{x})$



Distribution shifts are classically broken down to three types


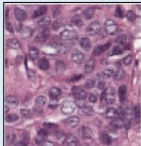
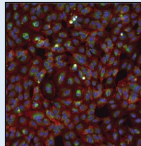

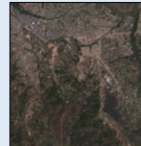
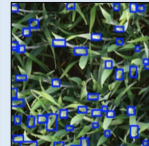
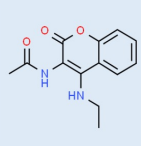

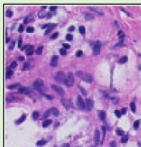
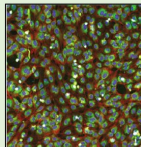



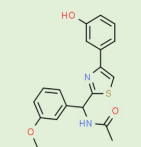
- In a supervised regime, we can write the joint distribution of data and labels as:

$$P(\mathbf{x}, y) = P(\mathbf{x}|y)P(y) \quad \text{-or-} \quad P(\mathbf{x}, y) = P(y|\mathbf{x})P(\mathbf{x})$$

- Covariate Shift: $P_{test}(y|\mathbf{x}) = P_{train}(y|\mathbf{x})$, but $P_{test}(\mathbf{x}) \neq P_{train}(\mathbf{x})$
 - Ex: $P_{test}(\mathbf{x})$ has more people over 60, but the per-person probability of polio has not changed
- Label Shift: $P_{test}(\mathbf{x}|y) = P_{train}(\mathbf{x}|y)$, but $P_{test}(y) \neq P_{train}(y)$
 - Ex: Everyone in P_{test} has been vaccinated. So, similar people still get polio, but it is less frequent
- Concept Drift: $P_{test}(y) = P_{train}(y)$, but $P_{test}(\mathbf{x}|y) \neq P_{train}(\mathbf{x}|y)$
 - Ex: Polio has mutated in P_{tgt} to affect younger instead of older people, but the *total* risk is the same

Distribution shifts are ubiquitous

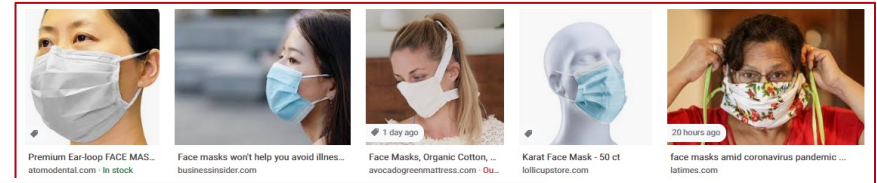
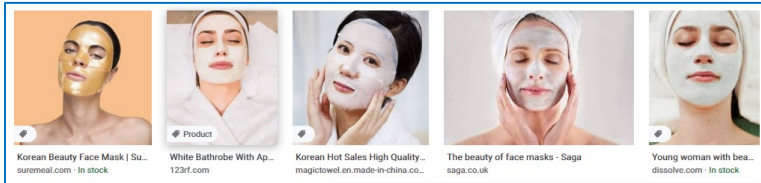
- Any changes in a current data generating environment can cause shifts
- Applying a model to a new domain is almost always a shift

Dataset	iWildCam	Camelyon17	RxRx1	FMoW	PovertyMap	GlobalWheat	OGB-MolPCBA	CivilComments	Amazon	Py150
Input (x)	camera trap photo	tissue slide	cell image	satellite image	satellite image	wheat image	molecular graph	online comment	product review	code
Prediction (y)	animal species	tumor	perturbed gene	land use	asset wealth	wheat head bbox	bioassays	toxicity	sentiment	autocomplete
Domain (d)	camera	hospital	batch	time, region	country, ru/ur	location, time	scaffold	demographic	user	git repo
Source example								What do Black and LGBT people have to do with bicycle licensing?	Overall a solid package that has a good quality of construction for the price.	<pre>import numpy as np ... norm=np.____</pre>
Target example								As a Christian, I will not be patronizing any of those businesses.	I *loved* my French press, it's so perfect and came with all this fun stuff!	<pre>import subprocess as sp p=sp.Popen() stdout=p.____</pre>

Exemplar Real-World Distribution Shift datasets from Stanford WILDS datasets overview [1]

Example: Google Search Results

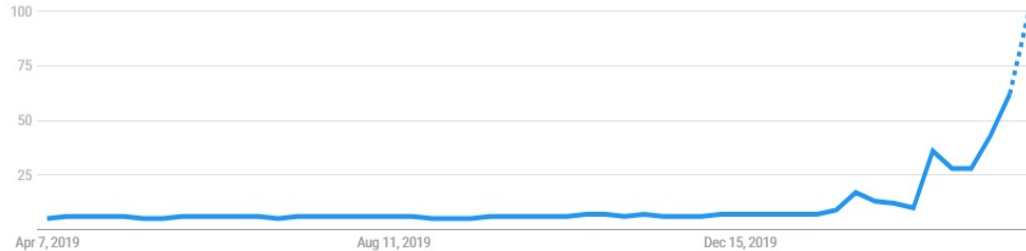
Face mask vs. Face mask?



Interest over time

Google Trends

● face mask



United States. Past 12 months. Web Search.

Part 2:

Living With Distribution Shifts

Detecting the problem and trying to answer *what happened?*

Detecting distribution shifts – common methods

- Distribution Shift detection answer the binary question: “Has a shift occurred?”
- Detecting distribution shift is a well-studied topic [3], most methods involve either:
 1. Statistical Hypothesis testing between P_{src} and P_{tgt} :
 $\phi(\hat{P}_{src}, \hat{P}_{tgt}) \geq \epsilon$, $\phi :=$ statistical divergence function (e.g., KL-divergence) and $\hat{P} :=$ a density model of the data (e.g., a normalizing flow)
 2. Training a domain classifier model f to classify between x_{src} and \hat{x}_{tgt} :
 $\mathbb{E}_{x \sim P_{tgt}}[f(x)] \geq \epsilon$, $\hat{x}_{tgt} :=$ an *estimate* of what samples from P_{tgt} will look like

We can use feature shift detection to localize the problem to specific features

- To detect feature shift [4], we define a conditional distribution hypothesis test:

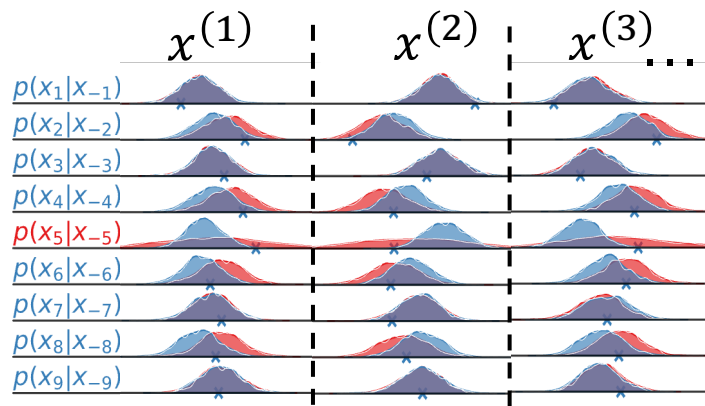
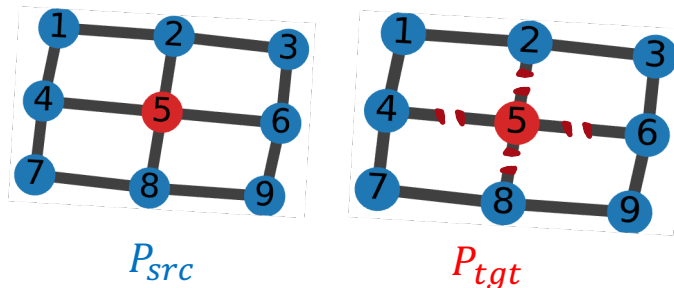
- $H_0: \forall \mathbf{x}_{-j} \in \mathcal{X}_{-j}, \hat{P}_{src}(\mathbf{x}_j | \mathbf{x}_{-j}) = \hat{P}_{tgt}(\mathbf{x}_j | \mathbf{x}_{-j})$

- $H_A: \exists \mathbf{x}_{-j} \in \mathcal{X}_{-j}, \hat{P}_{src}(\mathbf{x}_j | \mathbf{x}_{-j}) \neq \hat{P}_{tgt}(\mathbf{x}_j | \mathbf{x}_{-j})$

- Feature shift can happen in two stages:

- Detection: Do the conditional distributions of \hat{P}_{tgt} differ from the conditional distribution \hat{P}_{src} ?
 - Localization: Which feature(s) have caused this difference?

Feature Shift Toy Problem



$P(\mathbf{x}_j | \mathbf{x}_{-j})$ conditional distributions for $P \in \hat{P}_{src}, \hat{P}_{tgt}$

Feature Shift Detection is fast with Fisher divergence

- Fisher divergence test statistic based on the score function, $\psi := \nabla_x \log(p(x))$

$$\phi_{Fisher}(p, q) \triangleq \mathbb{E}_{p(x)+q(x)} [\|\psi(x; p) - \psi(x; q)\|^2] = \mathbb{E}_{p(x)+q(x)} \left[\left\| \nabla_x \log \frac{p(x)}{q(x)} \right\|^2 \right]$$

- Can compute multiple feature test statistics simultaneously

- $\phi_{Fisher}(p_{x_j|x_{-j}}, q_{x_j|x_{-j}}) = \mathbb{E}_{p(x)+q(x)} [(\psi(x; p) - \psi(x; q))^2]_j$

- Only a **single forward and backward pass** is needed to compute all conditional score functions, which is already done when updating a density model, \hat{P}

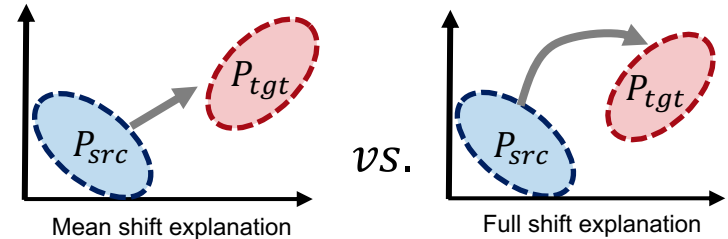
- Feature Shift tells us: “Has a shift occurred?” + “What set of features shifted?”

A distribution shift has been detected...now what?

We need to know more to respond effectively

- Problem: Once a shift has been detected, an operator needs to figure out what has changed in order to effectively respond
- Current simple approach: See how the means have shifted, $\mu_{src} - \mu_{tgt}$

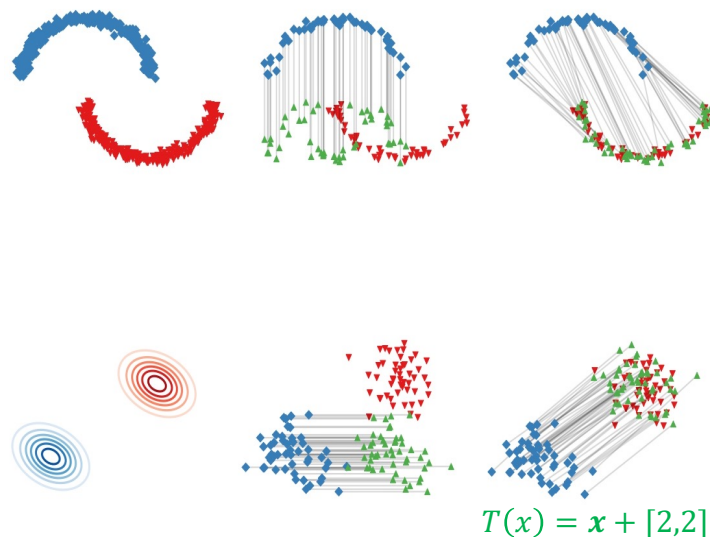
- Gives a rough approximation of shift
- However, this can miss important information:



- Our goal: Aid the operator by explaining how P_{src} shifted to P_{tgt}

Distribution shifts can be explained by hypothesizing how to map P_{src} to P_{tgt}

- Given two distributions P_{src} , P_{tgt} :
 - a transport map $T(\cdot)$, is a function which moves a point from P_{src} to P_{tgt} , such that $T_{\#}P_{tgt} \approx P_{src}$
- If T is interpretable, it can explain how P_{src} shifted to P_{tgt}



We can leverage prior Optimal Transport work to find **good** interpretable mappings

- Optimal Transport finds a minimum cost mapping T_c that aligns two distributions [10]
- By relaxing alignment and restricting our possible mappings to be interpretable we get *intrinsically interpretable transport* T_{IIT} :

$$T_{IIT} := \operatorname{argmin}_{T \in \Omega_{int}} \mathbb{E}_{P_{train}} [c(\mathbf{x}, T(\mathbf{x}))] + \lambda \phi(P_{T(\mathbf{x})}, P_{test})$$

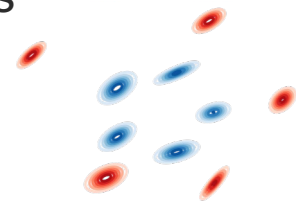
where $\Omega_{int} = \{T: \text{ s. t. } T \text{ is interpretable}\}$, $c(\cdot, \cdot)$ is a cost function (e.g., ℓ_2), and ϕ is a divergence

- T_{IIT} gives us a mapping which is faithful ($P_{T(\mathbf{x})} \approx P_{test}$), interpretable ($T \in \Omega_{int}$), and simple (minimizes a transport cost)
- Ω_{int} can be defined based on context, or one can use our pre-defined sparse-feature mappings or cluster-based mappings [5]

T_{IIT} can be used to gain actionable insights from explanations of complex shifts

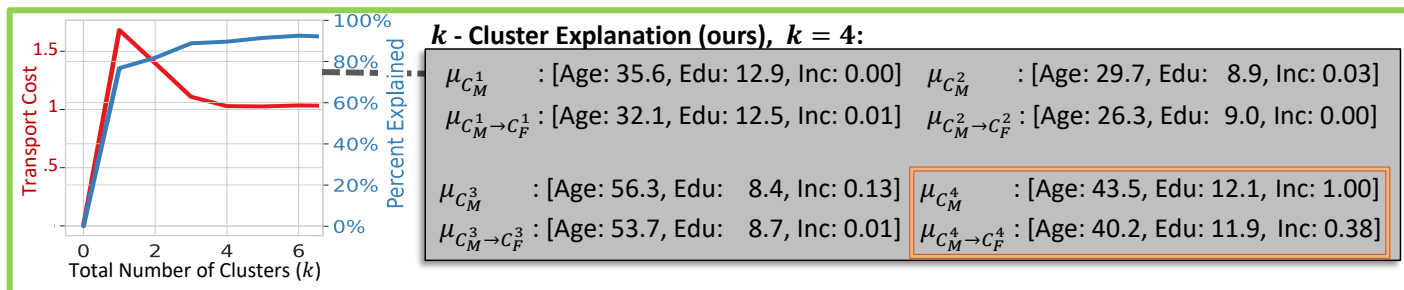
- Using our k -cluster mappings $\Omega_{cluster}^k$, we can see how heterogeneous groups (clusters) moved differently under a distribution shift

$$\Omega_{cluster}^k = \{T: T(\mathbf{x}) = \mathbf{x} + [\Delta]_c\}, \text{ where } \Delta \in \mathbb{R}^{dxk}, c = [k]$$



Example of a distribution shift

- We can use $\Omega_{cluster}^k$ to compare male and female responses to the 1994 US Census



Insight 1:
Income is largest predictor between M and F

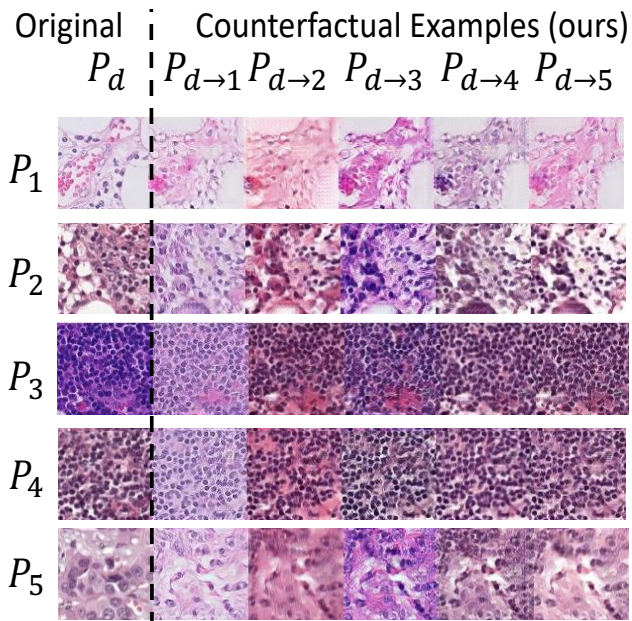
Insight 2:
The income difference is largest in M_{C^4} , middle-aged adults with a bachelor's degree

Transport Maps can also explain distribution shifts in high-dimensional regimes (images)

- When raw features are not semantically meaningful, but samples are (e.g., images), we can use post-hoc methods to understand T such as:

Distributional-Counterfactuals := $\{x, T(x) : x \sim P_{src}, T(x) \sim P_{tgt}\}$

- We can use distributional-counterfactuals to explain how H&E staining of tissue samples change across multiple hospitals [6]



Using StarGAN [7] to show the difference between tissue samples across 5 hospitals

Take-Aways on Distribution Shifts

- Distribution Shifts are ubiquitous, complex, and problematic for ML
- To combat distribution shifts we need to:
 1. Detect a shift has happened
 - Perform statistical hypothesis testing between \hat{P}_{src} and \hat{P}_{tgt} e.g., check for feature-shift
 2. Understand what the distribution shift has changed
 - Solve for a distribution shift explanation T and see if the changes are problematic
 3. React to the fix the shift
 - Possibly retrain models, fix the change in our environment, update training set, etc.

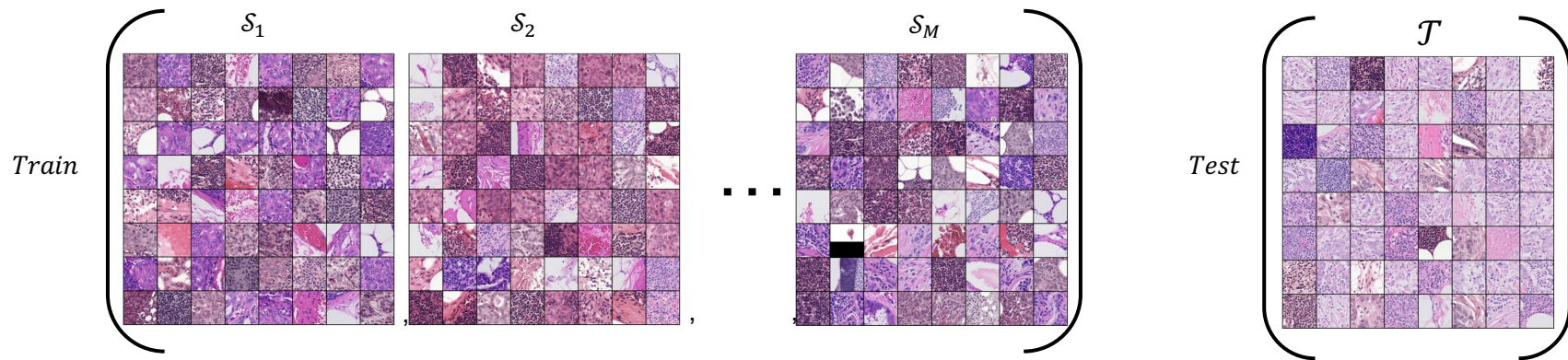
Part 3:

How to Avoid Problems with Distribution Shifts

Turning the problem into the solution – methods for domain generalization.

We can use sets of shifted distributions to build robust models

Given: M training domains $\mathcal{S} = \{\mathcal{S}_i \mid i = 1, \dots, M\}$ where $\mathcal{S}_i = \{(x_j^i, y_j^i, i)\}_{j=1}^{n_i}$



Goal:

- Find a model which can achieve a minimum error on an **unseen** test domain,

$$\mathcal{T} = \{x_j, y_j\}_{j=1}^{n_t}$$

- $\min_h \mathbb{E}_{(x,y) \in \mathcal{S}_{test}} [\ell(h(x), y)]$ for some loss function $\ell(\cdot)$ and $P_{XY}^{test} \neq P_{XI}^i$

Taxonomy of DG

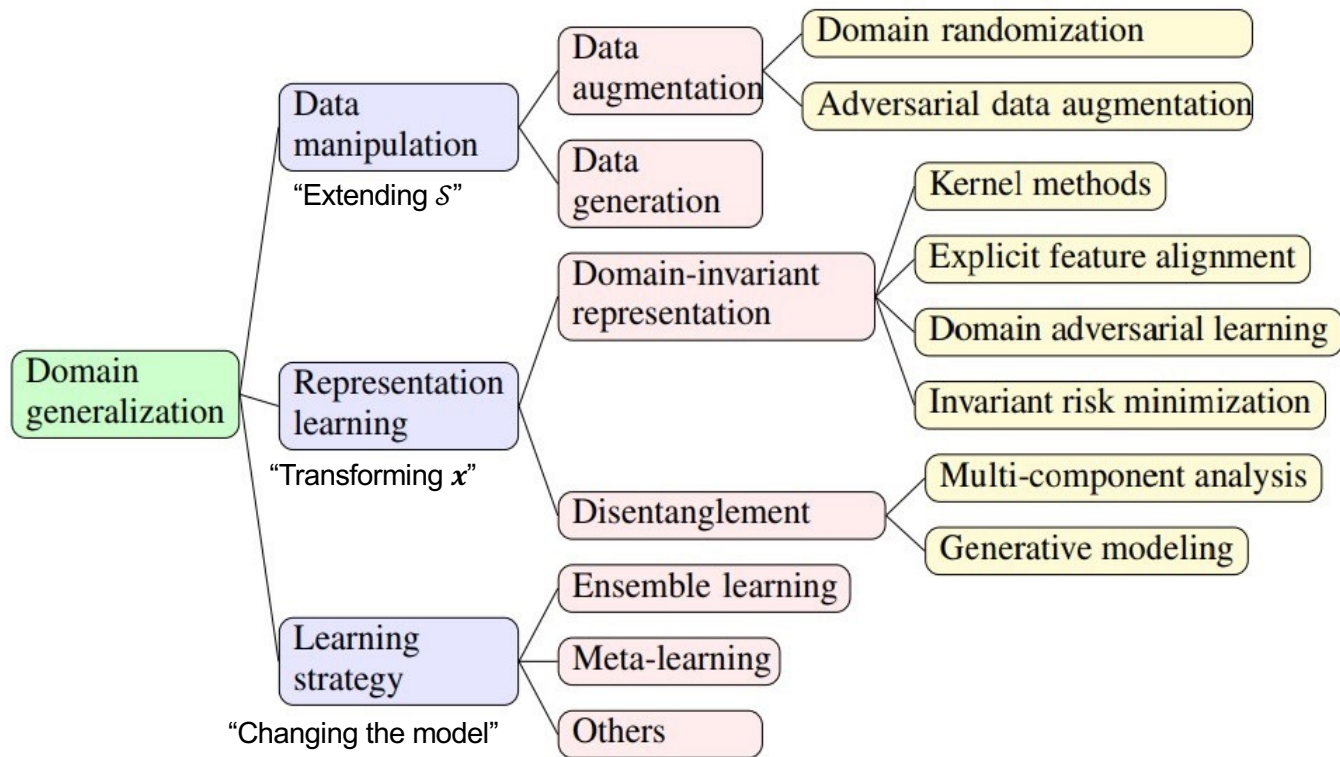


Image from [8].

An optimal representation which is invariant across domains in \mathcal{S} should generalize to unseen domains

- Domain adversarial learning

- Adversarial optimization where d discriminates the original domain of $g(x)$, and g finds a representation which aids the classifier $f(g(x))$ while fooling the discriminator

- $$\operatorname{argmin}_{f,g} \operatorname{argmax}_d \sum_{j=1}^M \sum_{(x,y) \in \mathcal{S}_j} \mathcal{L}_{f,g}(f(g(x)), y) + \mathcal{L}_d(d(g(x), j))$$

- Explicit feature alignment

- Alignment of the domain distributions using a shared feature extractor g

- $$\operatorname{argmin}_{f,g} \sum_{i \neq j}^M \operatorname{dist}(g_{\#}(\mathcal{S}_i), g_{\#}(\mathcal{S}_j)) + \sum_{(x,y) \in \mathcal{S}_i} \mathcal{L}(f(g(x)), y),$$
 where $\operatorname{dist}(\cdot, \cdot)$ is some notion of a distance or statistical divergence metric

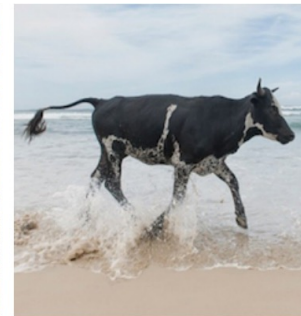
- Common representation functions: kernel methods, batch-instance normalization, neural networks

- Invariant risk minimization

- Find a data representation such that the optimal classifier $f^*(g(x))$ is the same across all environments



(A) Cow: **0.99**, Pasture: 0.99, Grass: 0.99, No Person: 0.98, Mammal: 0.98

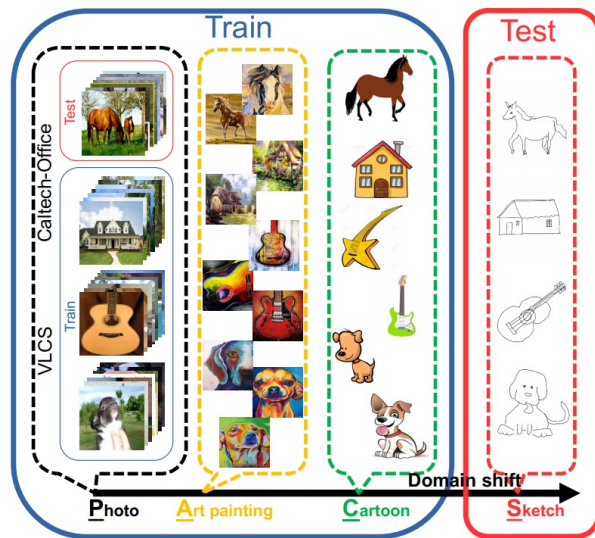


(B) No Person: 0.99, Water: 0.98, Beach: 0.97, Outdoors: 0.97, Seashore: 0.97

From: [Recognition in Terra Incognita](#) [10]

Feature-disentanglement learns both domain specific and domain-invariant representations

- Goal: learn function(s) that decompose samples into meaningful domain invariant $g_i(x)$ and domain specific features $g_s(x)$
 - $$\operatorname{argmin}_{g_s, g_{if}} \mathbb{E}_{x, y} \mathcal{L}(f(g_s(x)), y) + \lambda \mathcal{L}_{recon}([g_s(x), g_i(x)], x) + \lambda \mathcal{L}_{reg}(g_s(x), g_i(x))$$
- Multi-component analysis
 - During training, learn a universal model $\theta^{(0)}$ and domain-specific models $\{\theta^{(j)}\}_{j=1}^M$, and for inference use functional combination of the two
 - UndoBias: SVM where $\mathbf{w}(x) = \mathbf{w}^{(0)}(x) + \mathbf{w}^{(j)}(x)$ where $j \in \{1, \dots, M\}$ and is found via $\hat{j} = d(x)$, where d finds the domain which x is most likely to have come from
- Generative modeling
 - Use VAEs to find a latent space with disentangled representations of domain information, category information, and other information



From: Deeper, Broader and Artier Domain Generalization

References

- [1] Koh, Pang Wei, et al. "Wilds: A benchmark of in-the-wild distribution shifts." *International Conference on Machine Learning*. PMLR, 2021.
- [2] Chan, Stanley. "Chapter 4, Learning Theory." *ECE 595 / Stat 598: Machine Learning*, <https://engineering.purdue.edu/ChanGroup/ECE595/files/chapter4.pdf>.
- [3] Rabanser, Stephan, Stephan Günnemann, and Zachary Lipton. "Failing loudly: An empirical study of methods for detecting dataset shift." *Advances in Neural Information Processing Systems* 32 (2019).
- [4] Kulinski, Sean, Saurabh Bagchi, and David I. Inouye. "Feature shift detection: Localizing which features have shifted via conditional distribution tests." *Advances in Neural Information Processing Systems* 33 (2020): 19523-19533.
- [5] Kulinski, Sean, and David I. Inouye. "Towards Explaining Distribution Shifts." *arXiv preprint arXiv:2210.10275* (2022).
- [6] Kulinski, Sean, and David I. Inouye. "Towards Explaining Image-Based Distribution Shifts." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022.
- [7] Choi, Yunjey, et al. "Stargan: Unified generative adversarial networks for multi-domain image-to-image translation." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018.
- [8] Wang, Jindong, et al. "Generalizing to Unseen Domains: A Survey on Domain Generalization." *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, International Joint Conferences on Artificial Intelligence Organization*, 2021, pp. 4627–35.
- [9] Beery, Sara, Grant Van Horn, and Pietro Perona. "Recognition in terra incognita." *Proceedings of the European conference on computer vision (ECCV)*. 2018.
- [10] Peyré, Gabriel, and Marco Cuturi. "Computational optimal transport: With applications to data science." *Foundations and Trends® in Machine Learning* 11.5-6 (2019): 355-607.

Thanks for listening :)

I'm happy to answer any questions you have now.

If you would prefer to chat after, just email me at: skulinski@purdue.edu , or
you may find answers/more ways to reach me on my website: seankulinski.com