

Towards Explaining Distribution Shifts

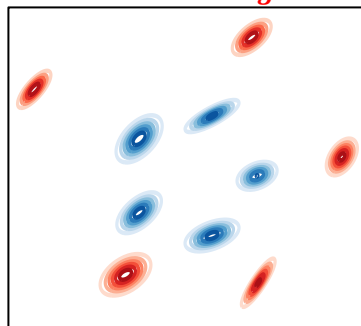


Sean Kulinski

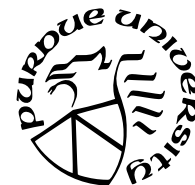
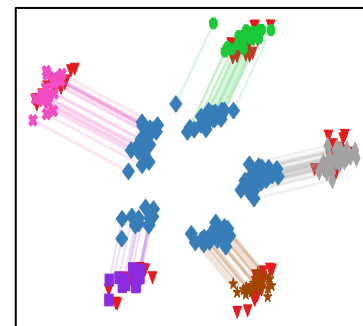


David Inouye

Oracle Shift from
 P_{src} to P_{tgt}



Proposed Distribution
Shift Explanation

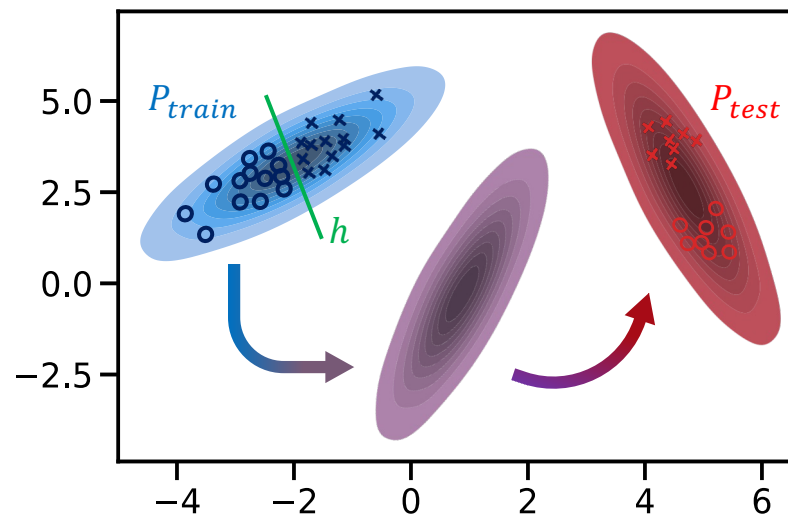


A distribution shift is when a data distribution changes from what is expected

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


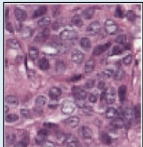
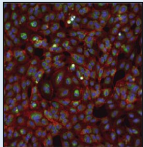



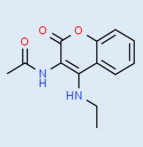

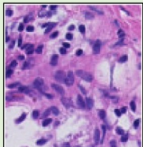
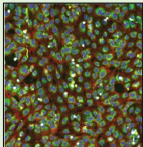



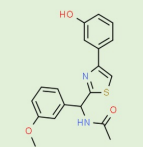
$$P_{test}(x) \neq P_{train}(x)$$

- Under distribution shift, the patterns learned by **a model** might not be present in P_{test}



Distribution shifts are ubiquitous

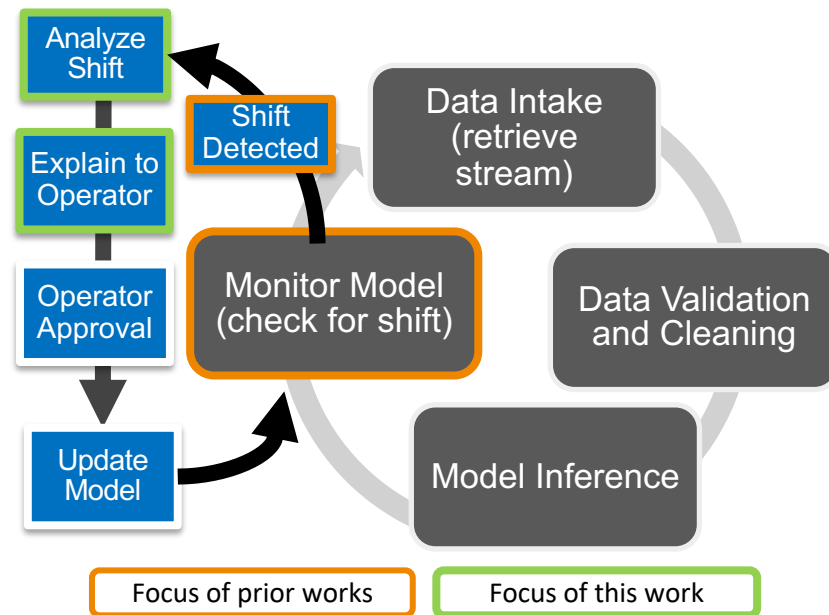
- Any changes in a current data generating environment can cause shifts
- Applying a model to a new domain is almost always a shift

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 benchmarks overview

Knowing what has changed under a shift allows us to more **effectively** respond to mitigate the shift

- **Problem:** Most prior works focus on only *detecting* a shift and do not help with “How should I respond?”
- To most effectively mitigate the shift, an operator needs to know what changed
 - E.g, “Preferences of 18-25 year-olds changed” or “X feature of the data intake pipeline is broken”
- **Our goal:** Aid the operator by **explaining** how P_{src} shifted to P_{tgt}



A typical ML deployment cycle

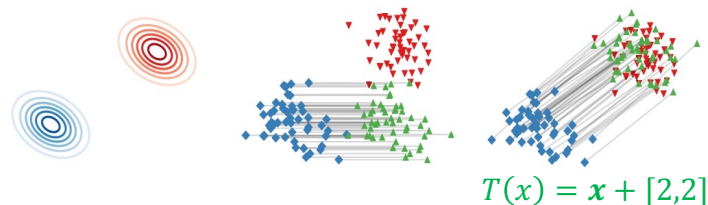
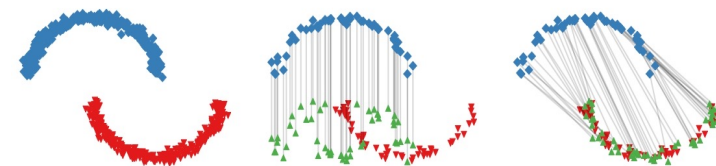
Distribution shifts can be explained by hypothesizing

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

$$P_{T(P_{src})} \approx P_{tgt}$$

- If T is interpretable, it can explain how P_{src} shifted to P_{tgt}



$$T(x) = x + [2,2]$$

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

- By relaxing alignment in Optimal Transport and restricting our possible mappings to be interpretable we get *Intrinsically Interpretable Transport*:

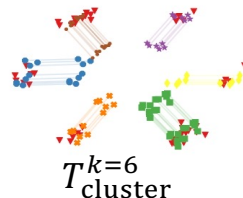
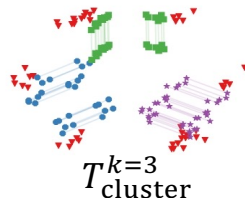
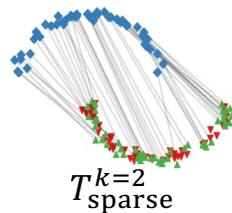
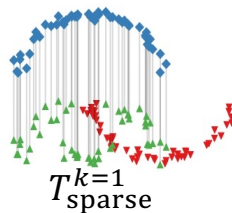
$$T_{IIT} := \operatorname{argmin}_{T \in \Omega_{int}} \mathbb{E}_{P_{train}} [c(x, T(x))] + \lambda \phi(P_{T(x)}, P_{test})$$

Ω_{int} : A set of interpretable mappings

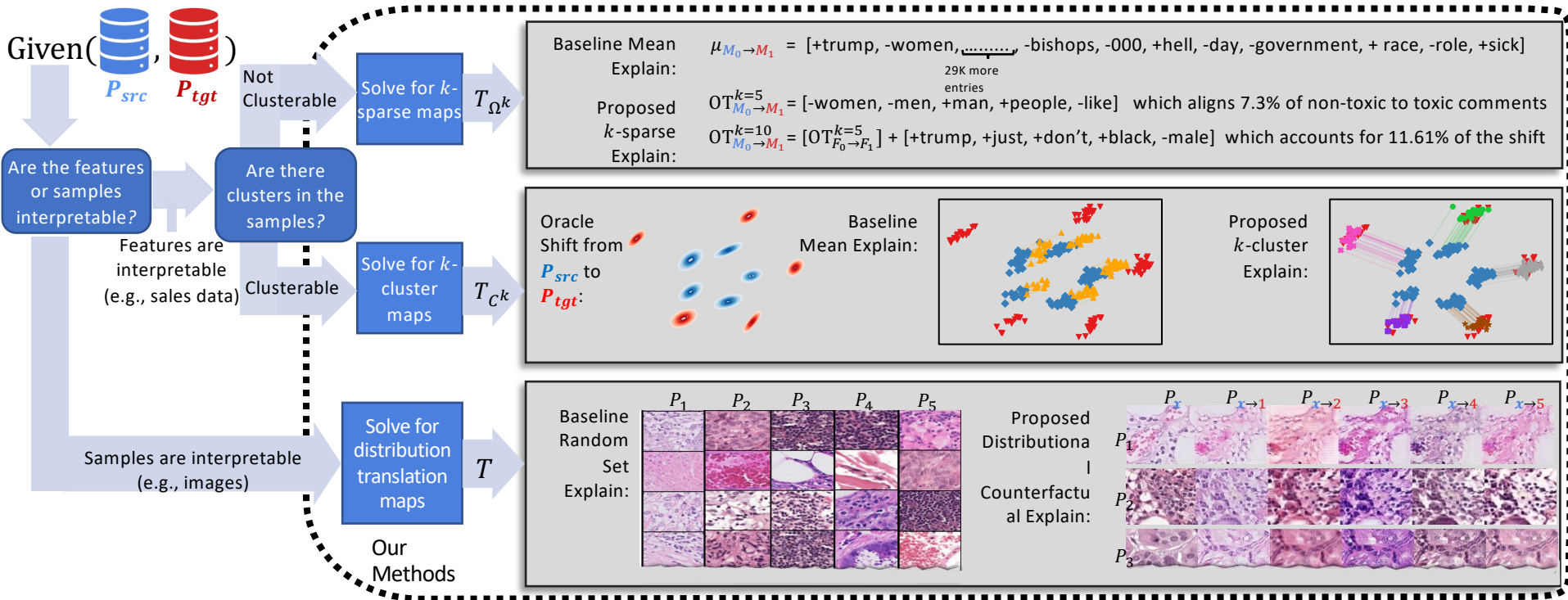
Cost function: T should retain as much of the original point as possible

Divergence: T should align $P_{T(x)}$ and P_{test} as much as possible

- Ω_{int} can be defined based on context, or one can use our pre-defined mappings: k -sparse feature mappings or k -cluster mappings



Methodology for solving for a shift explanation

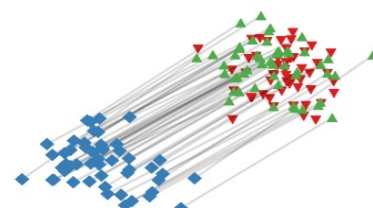
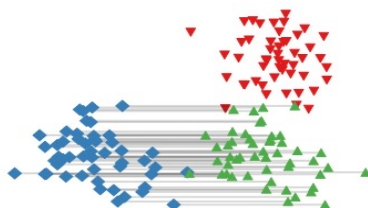
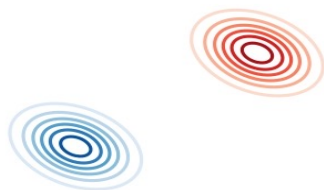


k -Sparse Feature Mappings can show how features moved along defined axes

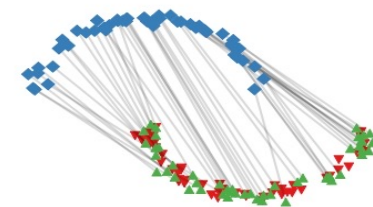
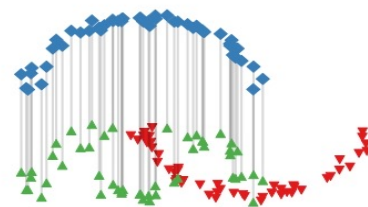
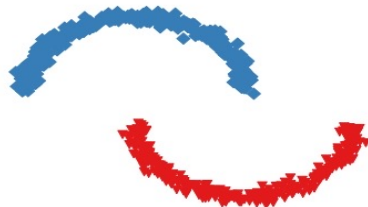
- Ω_{sparse}^k : Find a T which yields the best alignment, while only moving points

along k dimensions

Simple mean shift:



Complex conditional shift:



Oracle Shifts from P_{src} to P_{tgt}

$T_{sparse}^{k=1}$

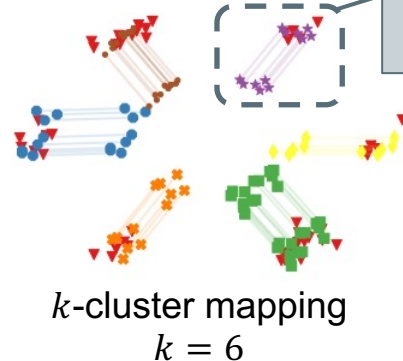
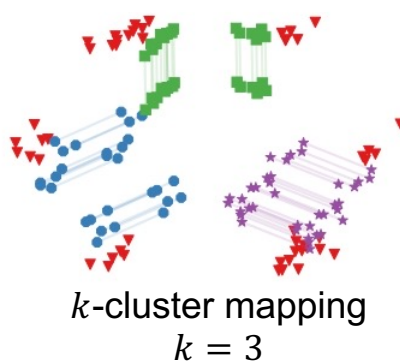
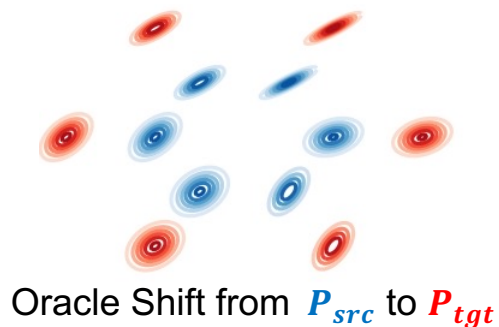
$T_{sparse}^{k=d}$

k -Cluster Mappings can show how heterogeneous subgroups have shifted

- $\Omega_{\text{cluster}}^k$: Find k -cluster-specific transport maps which maximizes alignment

between $P_{T(P_{tgt})}$ and P_{tgt}

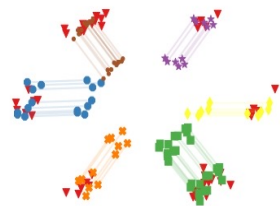
- We can restrict per cluster transport maps to a specific class of transport functions



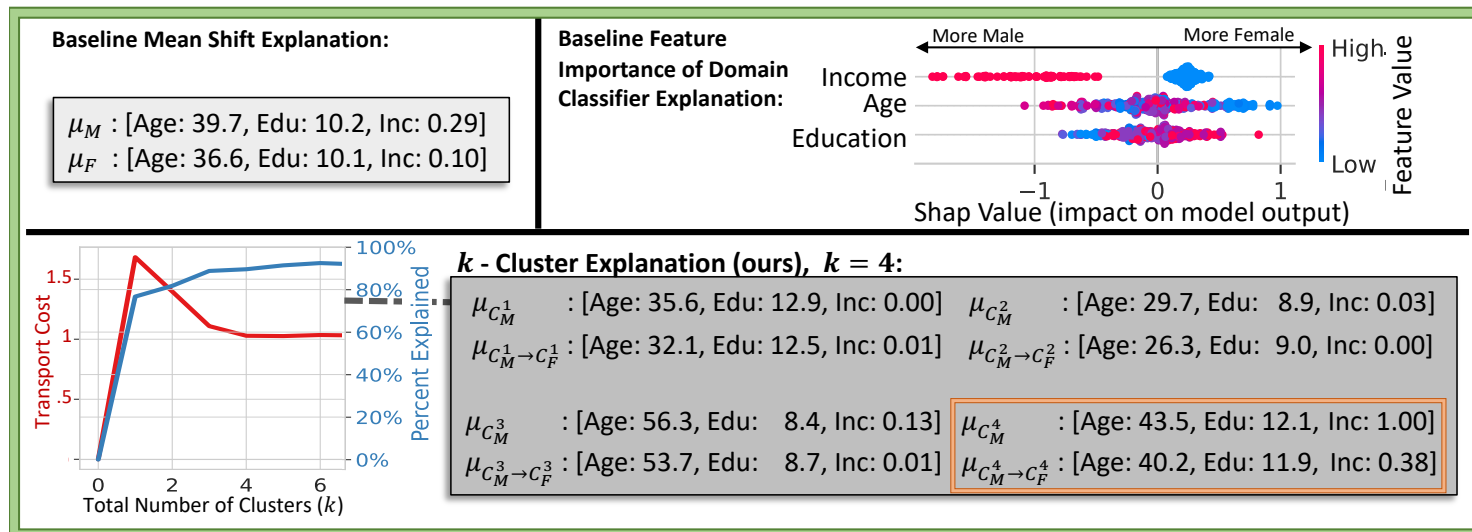
Moving points by cluster-specific vector:
 $T^{(C)} = x + \delta^{(C)}$
if $x \in C$

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 heterogenous groups (clusters) moved differently under a distribution shift

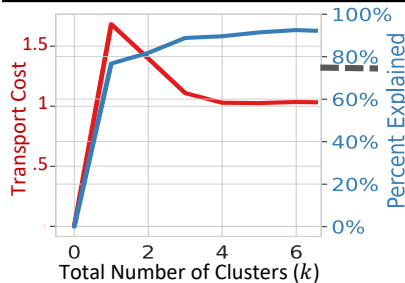


Example of 6-cluster mapping



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



Using $\Omega_{cluster}^k$ to compare male and female response to the US 1994 Census

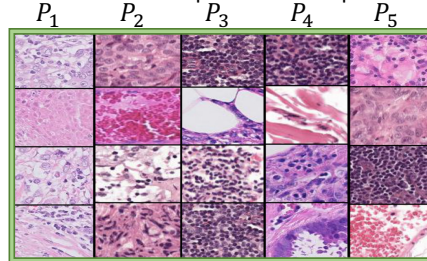
Transport Maps can also explain distribution shifts in high-dimensional regimes (images)

- When raw features are not semantically meaningful, but samples are (e.g., images), we can use *domain counterfactuals* to understand a complicated T

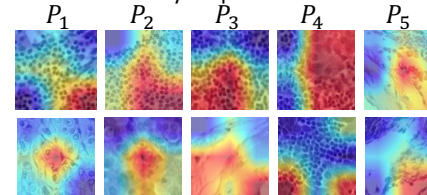
- Distributional-Counterfactuals :=

$$\{x, T(x): x \sim P_{src}, T(x) \sim P_{tgt}\}$$

Baseline: Visual Inspection of Samples



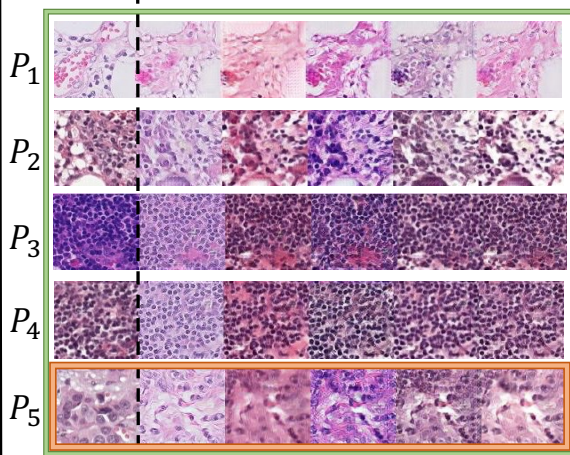
Baseline: Saliency Maps for Domain Classifier



Insight ①:

There seems to be a difference in staining across hospitals

Original | Counterfactual Examples (ours)
 P_d | $P_{d \rightarrow 1}$ $P_{d \rightarrow 2}$ $P_{d \rightarrow 3}$ $P_{d \rightarrow 4}$ $P_{d \rightarrow 5}$

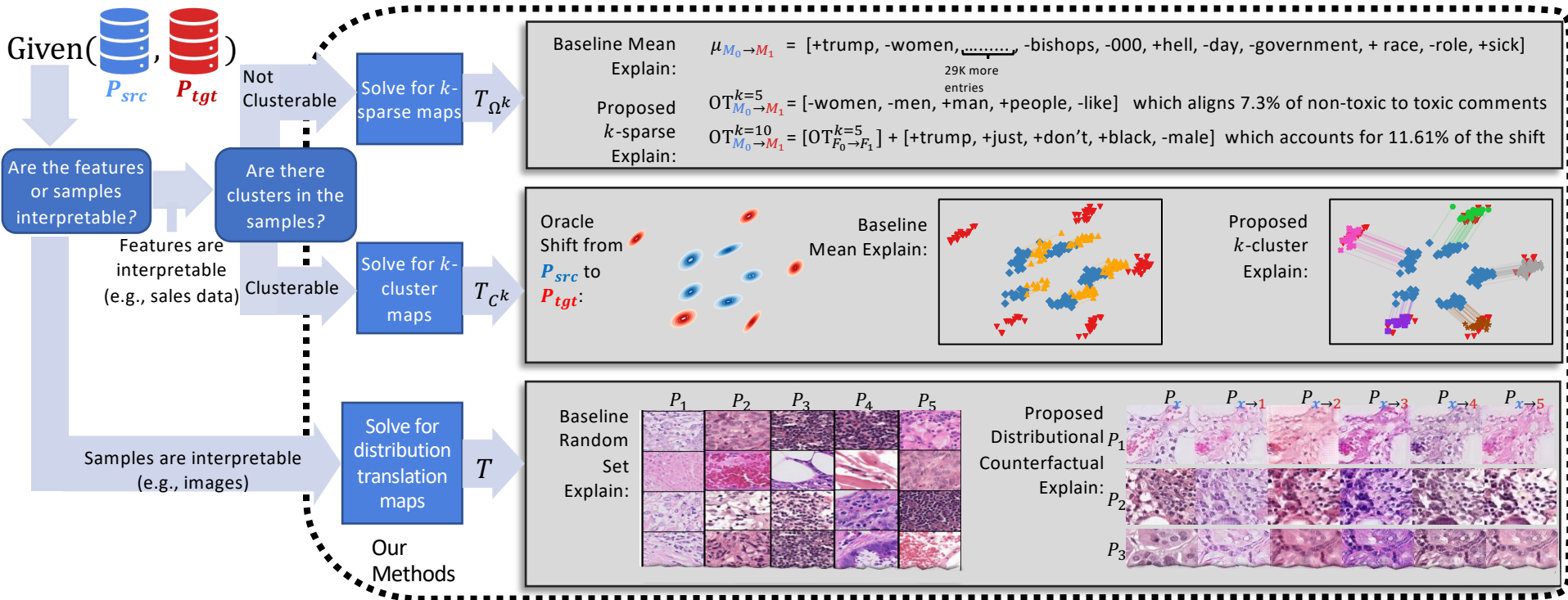


Insight ②:

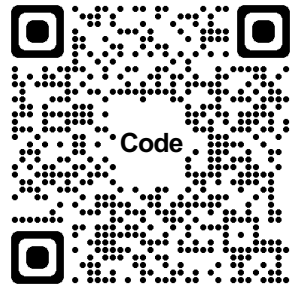
There is a clear difference in staining, and it seems to be unique to each hospital

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

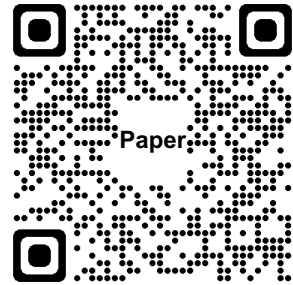
Methodology for solving for a shift explanation



Thank you for listening!



Towards Explaining Distribution Shifts



Sean Kulinski



David Inouye

