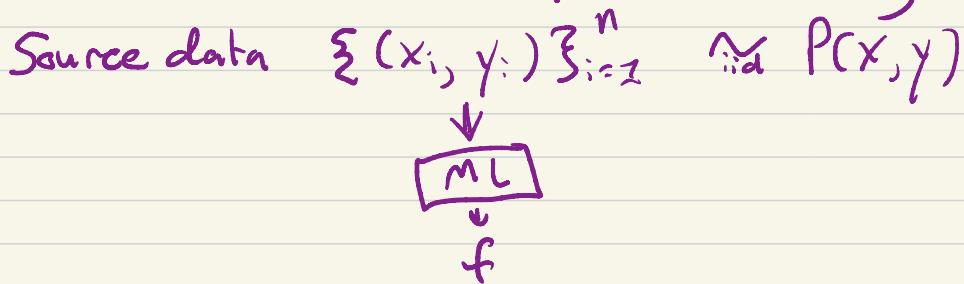


10732 Structured Approaches to distribution shift + the label shift story

Refresher: standard prediction setting



Deployment/target data $\{(x_j, y_j)\}_{j=1}^m \sim_{\text{iid}} P(x, y)$
(observe only $\{x_j\}_1^m$)

Because $P_S = P_T$, licensed to
evaluate on holdout set of source data
(exchangeable w deployment data).

Here, fairness of model on target dist is
directly observed, no identification problems.
(licensed to try (almost) anything, see what works)

What can we do when $P_{src} \neq P_{tgt}$?

Some lenses:

- theoretical models
 - distribution robustness over uncertainty sets
 - identification under structure
- Deep learning playbook
 - benchmarks + tricks
- Look out at the real world

Themes:

- All assumptions are too strong
- Problems ill-posed absent assumptions
- Lack a distribution over distributions
- Effective sample size
- What role do theory, benchmarks, etc. play?

Structure of a Distribution Shift Problem:

- Set of environments / domains
- Underlying structure shared across environments
- Observability - what variables seen where?
- Manipulation rules:
 - what can change?
 - in what way?
 - by how much?
 - in vs. out-of-support?
- Statistical / inferential capabilities
 - what relationships can be learned?
 - how effectively?
 - by what measure?
- Objectives(s)

Steps:

- Identification
- Estimation / Learning

Note - we are interested in
"practical" identification strategies
involving only those quantities we
have the power to estimate,
approximated to the degrees
actually attainable

Identification Tips:

- Worry about identification before worrying about practical procedures
- First focus on finite data domains w/ discrete/tabular data,
- Assume direct access to joint dist over observables
(If we can't handle this, we're in trouble)

The Move: Leverage blackbox predictors as nonparametric hummers

- Fashion an "iid problem" for DL, where fitness can be assessed on appropriate holdout data
- Show that if net hits some reasonable property, then we can work in push-forward space.
- Use net + structure to identify/estimate target parameters.

Label Shift



Source: x, y visible

Target: x visible

- Motivations:
- disease propagation
 - general class balance problems

Side note: Test item effect (Zhu et al 2010)

Formally

- $P(X|y) = Q(X|Y)$
- $\text{Supp}(Q(y)) \subseteq \text{Supp}(P(y))$

Identification:

Everything we see:

$$P(x|y=1) \dots P(x|y=k) = Q(x)$$

Linear system, solution identified whenever class-conditionals are linearly independent.

$$Q(x) = \sum_y Q(x,y) = \sum_y Q(y) P(x|y) = \sum_y Q(y) P(x|Y)$$

BBSE:

- Train neural net f on source data
- Consider mean classifier output

$$\begin{matrix} C_y^1 \\ \hat{y} \end{matrix} = \begin{matrix} Q(y) \\ M_f^Q \end{matrix}$$

Detection: $f(x)$ is sufficient to detect
label shift

$$\text{Estimation: } \hat{\theta}(y) = \hat{C}_{\hat{y}, y} M_f^Q$$

$$\hat{w}(y) = \hat{C}_{\hat{y}, y} M_f^Q$$

Error bounds: w.h.p, for some constant C ,

$$\| \hat{w} - w \|_2^2 \leq \frac{C}{\sigma_{\min}^2} \left(\frac{\| w \| \log n}{n} + \frac{k \log m}{m} \right)$$

(SAERENS 2006)

Other estimation ideas using Black Box:

EM / MLLS

1. Start off as though $\hat{P}_t(Y) = P_S(Y)$

For $n=1 \dots$

$$\hat{P}_t^n(Y_i | X_i) = \frac{\hat{P}_S(Y_i | X_i) \hat{P}_t^n(Y_i)}{\sum_k \hat{P}_S(Y_k | X_i) \frac{\hat{P}_t^n(Y_k)}{P_S(Y_k)}}$$

$$\hat{P}_t^{n+1} = \frac{1}{N} \sum_{i=1}^n \hat{P}_t^n(Y_i | X_i)$$

• Exactly the MLE of $P_t(Y)$ when $\hat{P}_S(Y | X)$ is the true source conditional,

Else, when is this consistent?

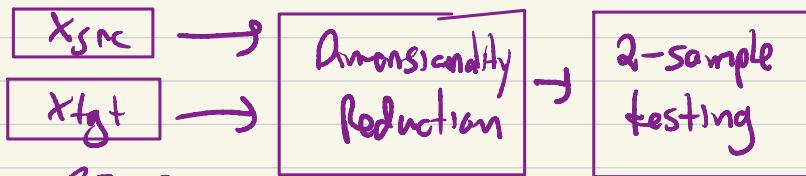
Turns out, requires some invertible CM
+ canonical calibration (Garg 2020)

In practice, works well in BCTs (^{Alexandri} 2019)
(temp scaling w/ per-class parameters)

Error decmp: $\|\hat{w}_p - w\| \leq \|w_p - \hat{w}\| + \|w_p - w^*\|$
finite samples calibration

(Robensler 2018)

Fairly loudly! pipeline for detecting shift



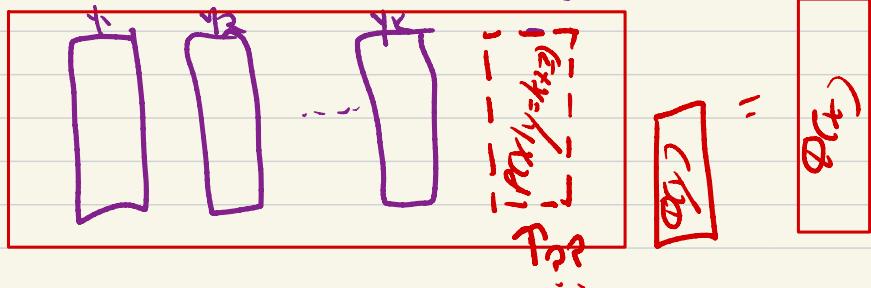
- + BBSD effective broadly, even w LS violated
- too strict, there's always a tiny shift, but when should we care?

IW-ERM \rightsquigarrow DL (Byrd 2019)

$$\min_{\theta} \sum_{i=1}^n w_i \ell(f_{\theta}(x_i), y_i)$$

Problem: what role do weights play when training loss for all examples $\rightarrow 0$?

Extending Label Shift to Handle a novel class



No source data for est. $P(x|y=k+1)$.
Who is to say $P(x|y=k+1) \neq Q(x)$?

Perhaps $Q(y=k+1) = 1!!$

Requires addtl identifiability constraints.

Focus on base case: $k=1$

$$\boxed{P(+|x_1)} = \boxed{(+)\theta}$$

don't need
to know which
 $x!$

Assumption: Separable Sub-Support

$$\exists x \in \mathcal{X} \quad P(x|y=+) > 0 \\ \quad \quad \quad \& \quad P(x|y=-) = 0$$

Steps:

① Mixture proportion estimation
of $Q(y=+) \hat{\Delta} \alpha$

② Learn a PvN classifier

Previous approaches

Elkan / Noto provide algorithms

key idea is to relate PvN quantities
to PvU "non-traditional classifier"

No theory / analysis.

Other recent approaches act in classifiers
output space but remain heuristic:
no guarantees for consistency.

Sugiyama introduced unlabeled loss
for learning the classifier:

$$\hat{R} = \alpha R_p^+ + R_u^- - \alpha R_p^-$$

Interpretation:

$$\alpha^* = \max_{\alpha} \alpha$$

$$\text{st } \forall x \quad \alpha P(x|y=+) \leq Q(x)$$

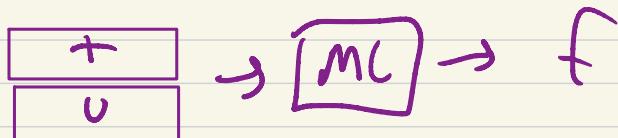
Intuition:

higher α : impossible

lower α : violates irreducibility / separable subspace

✗ Estimation Strategy:

Domain discrimination

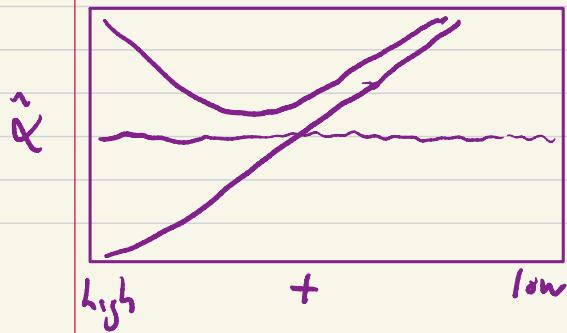


Insight: α determined by the "most positive" points (those in separable sub-support)

- Use f to map all points to scalar $z_i = f(x_i)$
 - Set threshold t
 - Using holdout + & 0 data estimate

$$\hat{\alpha} = \frac{P_0(z > t)}{P_+(z > t)} \geq \alpha$$

Bias - Variance tradeoff



Choose t that gives lowest upper bound on α

Bound

$$\|\hat{\alpha} - \alpha\| \leq \frac{C_1}{C_p C_2} \left(\sqrt{\frac{\log(t/8)}{n_q}} + \sqrt{\frac{n}{n_p}} \right)$$

Obtaining the optimal classifier:

A nPU loss

B NNnPu loss

Proposed CVIR (for separable data)

Sort data, discard highest loss & ^{fore-tick} among unlabeled at each iteration.

Consistent when data separable
avoids overfitting mislabeled data

Need more innovation for non-sep data

Further empirical improvements w/
iterated $\hat{\alpha}_i, \hat{f}$ estimation (TED_i^N)

Back to Open Set Label Shift (Garg 2022)

- Reduction to k PU problems
- Require each present class has a separable sub-support

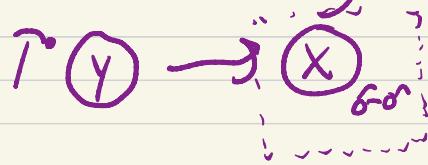
OR negatives entirely separated.

- Tricks for coping with overestimation of $\hat{\alpha}$

Algorithm :

- ① Trim source classifier f
- ② Estimate each $\hat{\alpha}_i$: $i \in \text{PU}$ via BBE, using $f(y=i|x)$ as score
- ③ Normalize to sum to 1 (correcting about for overestimation)
- ④ Re-weight source data via $\hat{\alpha}_i$
- ⑤ Estimate frac belonging to novel class through one final application of BBE, treating the reweighted src as one giant + class.

Unsupervised Learning under Latent Label Shift



Data observed across many domains
(# domains \geq # labels)

Given this structure (shift due to label shift),
when are underlying labels recoverable?
(up to permutation)

Motivation:

- Butterflies, caterpillars

Principle:

- items that shift together group together

Identification

- Finite data, isomorphism
to topic modeling, apply Arora

↑ Separable subgroups = anchor

- Continuous data + block box
→ product $\Phi(d|x)$