

# Research PR 1 - Counterfactual Explanations for ML October 2020

• link bet<sup>n</sup> wt cud happen given i/p changed in specific way  
 Post-hoc explanations for opaque models → model specific  
 → model agnostic

\* **QII** Quantative Input Influence  
 sees feat. correlations by randomly changing 1 feat at a time & calc. avg. marginal expected contribution of changed variable on model output over others

\* **SHAP** Shapley Additive explanations  
 game theoretic approach + local explanations  
 SHAP = f(LIME, Shapley sampling, QII, DeepLIFT, Layerwise Regr propaga<sup>n</sup>)

feature importance → SHAP  
 → QII  
 visual {  
 → Partial Dependency  
 → Accumulated Local Effects  
 → Individual Cond<sup>n</sup> Expectation  
 local approx {  
 → LIME  
 → Anchor  
 local example based → Counterfactual Exp

Kernel Gradient → TreeExplainer

\* **LIME** Local Interpretable Mod-agnostic Exp

local surrogate linear model on perturbed data  
 local fidelity ⇒ ↑ local fit (loss has complexity penalty)  
 typically uses exponential smoothening & lasso.

\* **PDP** - Expected target response wrt x keeping all else constant

\* **ALE** - small window, handles corr<sup>n</sup>

\* **Anchors** - Lime with if-else (nonlinearity to handle unseen test instances)  
 → Enemy vs friendly tank  
 → Husky vs Wolf (snow)

Others - Goldeneye, PALM, CAM, Grad-CAM, MES

\* Not driver analysis but wt Δinput; pushes record; across decision bound<sup>y</sup>

CF explanations

Similar terms -

- Validity set up as optim problem, dist based
- Actionability Only change continuous, bias in cat<sup>g</sup>
- Sparsity Penalty that ↑ sparsity in diff of CF & data
- Data Manifolds Penalty that ↑ adherence to manifold
- Causality Retain causal relations
- Amortized Inference opti<sup>x</sup> expensive ⇒ generative technique
- Speed for linear/piecewise can use LPP to create CFs

- Recourse (CFs better)
- Inverse classification ~
- Contrastive explanation !N
- Adversarial Learning (CFs more parsimonious Advs can hv ↑ confidence misclassification)

Validation - ① Validity ② Proximity ③ Sparsity ④ Diversity ⑤ Close to train ⑥ Causal IM1 and IM2

Proportion of CFs actually going across decision boundary

CF vs record Minkowski or Mahalanobis

feat. ↓ change

at least 1 CF is useful

avg dist to k neighbors local outlier factor or reconstruction error wrt VAE (train)

auto-encoders built class wise & check ~ reconstr<sup>n</sup> errors



Research Ppr 2 - Model agnostic Counterfactual Exp using RL June 2021  
 Replace expensive opti with end-to-end learnable RL  
 + Model Agnostic + target conditional CF instances + Num & lat while preserving immutability wrt actionability

Many Opti  $\rightarrow$  Fewer Opti  $\rightarrow$  Reconstruction  $\rightarrow$  utilize pre-trained  $\rightarrow$   
 (Batch level) losses via AEs GANs or VAEs to skip opti

- Other ways - LORE  $\rightarrow$  class cond<sup>n</sup> generative DICE
- Local Surrogate model around instance uses determinantal point processes
  - Train on syn. data created by using GAN  $\rightarrow$  expects diverse CFEs joint intensity
  - Decision rules from local won change model prediction leveraged for counterfactual Wud DICE be overfitted? equates to determinant of a function
- Even if u do class conditional generative model for batches of CFEs, this isn't feasible for non-differential models as backprop of gradients needed

RL-usage: Identify optimal  $\phi$  function (agent reward)  
 Deep Deterministic Policy Gradient bypasses computation by some approx. (OpenAI ppr)

Say  $x$  = instance  $x_{CF}$  depends upon  $(x, y_M, y_T, c) = s$  (state) &  $\mu(s)$  is actor network  
 $M$  = black box model  
 $y_M = M(x)$  prediction Reward classify = 1

$x_{CF} = x + \delta_{CF} \rightarrow$  sparse Reward regression  $\propto$  proximity  
 $\rightarrow$  counterfactual perturbation \* uses autoencoder help  
 instance \* perturbations in latent space  
 optimal conditioning vector  $c$  \*  $X$  shift & scale invariant then sampled from Beta(2,2) [feature conditioning++]  
 GA models  $x_{CF}$  allowing for  $c$  || categorical from Bern(0.5)  
 generative

Baselines for comparison - Model Agnostic (LORE)  
 \* DICE with random perturbations  $\Rightarrow$  better validity, smaller proximity  $L_0$ , comparable  $L_1$ , worse in-distributionness  
 ? Gradient-free (MO)  
 Tabular-oriented (DICE)  $\rightarrow$  random genetic

Pros: Cons:  
 Better in-distributionness (MM)  
 Better transfer, scale aspects  
 Lesser time compute  
 Easier feature conditioning  
 No need for smaller problems given  
 DICE is almost at-par