

Distribution Shift in Deployment Impacts Performance



Main Contribution

Detecting **performance drops on** new, unlabeled user data is a key challenge for deployed ML models. Our suitability filter compares model outputs from new and old data via statistical testing to determine if the model's accuracy has **degraded** beyond a user-defined performance margin. Our testing procedure even guarantees bounded **FPR** under some distribution shifts.



How Can a Model User Pick a Good Model?



Suitability Filter: A Statistical Framework for Classifier Evaluation in Real-World Deployment Settings

Angéline Pouget¹ Mohammad Yaghini² Stephan Rabanser² Nicolas Papernot² ¹ETH Zürich ²University of Toronto & Vector Institute

 $\frac{1}{2} \sum \mathbb{I}\{M(x) = \mathcal{O}(x)\} > 0$

$$\overline{D_u} \sum_{x \in D_u} \mathbb{I}\{M(x) = O(x)\} \ge \overline{|I|}$$

How can we proxy accuracy on the user data without access to labels?

A suitability filter $f_s : \mathcal{X} \to \{\text{SUITABLE}, \text{INCONCLUSIVE}\}$ outputs SUITABLE iff M is suitable for use on D_u with high probability and INCONCLUSIVE otherwise.



Suitability Signals



$$\sum \quad \mathbb{I}\{M(x) = y\} - m.$$

Margin Adjustment

$m' = m + \Delta_{\text{test}} - \Delta_{u}$ 0.9 D_{test} E 0.8 m' **상** 0.7 0.6 0.5 ${m}$ 0.7 0.8 0.4 0.5 0.6 Ground Truth Accuracy

Suitability Filters Work Across Domains

Dataset	Acc
FMoW ID	$81.8_{\pm 3.1}\%$
FMoW OOD	$91.9_{\pm 2.5}\%$
RxRx1 ID	$100.0_{\pm 0.0}\%$
RxRx1 OOD	$97.5_{\pm 7.2}\%$
CivilComm ID	$93.3_{\pm 5.3}\%$

On FMoW, we detect performance deterioration of > 3% with 100% accuracy.

