

# Pitfalls of evaluating a classifier's performance in high energy physics applications

Gilles Louppe, NYU ([@glouppe](#))  
Tim Head, EPFL ([@betatim](#))

December 11, 2015  
ALEPH workshop, NIPS

# Disclaimer

The following applies **only** for the learning protocol of the *Flavours of Physics* Kaggle challenge (Blake et al., 2015).

See (Louppe and Head, 2015) for the notebook explanations.

# Flavours of Physics: Finding $\tau \mapsto \mu\mu\mu$ challenge

Given a learning set  $\mathcal{L}$  of

- simulated signal events  $(\mathbf{x}, s)$
- real data background events  $(\mathbf{x}, b)$ ,

build a classifier  $\varphi : \mathcal{X} \mapsto \{s, b\}$  for distinguishing  $\tau \mapsto \mu\mu\mu$  signal events from background events.

$$\varphi^* = \arg \min_{\varphi} \frac{1}{N} \sum_{\mathbf{x}_i} L(\varphi(\mathbf{x}_i))$$

## Control channel test

The simulation is not perfect: simulated and real data events can often be distinguished.

To avoid exploiting simulation versus real data artefacts to classify signal from background events, we evaluate whether  $\varphi$  behaves differently on simulated signal and real data signal from a control channel  $\mathcal{C}$ .

Control channel test: Requires the Kolmogorov-Smirnov test statistic between  $\{\varphi(\mathbf{x})|\mathbf{x} \in \mathcal{C}^{\text{sim}}\}$  and  $\{\varphi(\mathbf{x})|\mathbf{x} \in \mathcal{C}^{\text{data}}\}$  to be strictly smaller than some pre-defined threshold  $t$ .

# Loophole

Assuming that

- control data can be distinguished from training data,
- simulated features are more discriminative than they are in real data,

Then, it exists classifiers  $\varphi$  exploiting simulation versus real data artefacts to classify signal from background events, for which **the control channel test succeeds**.

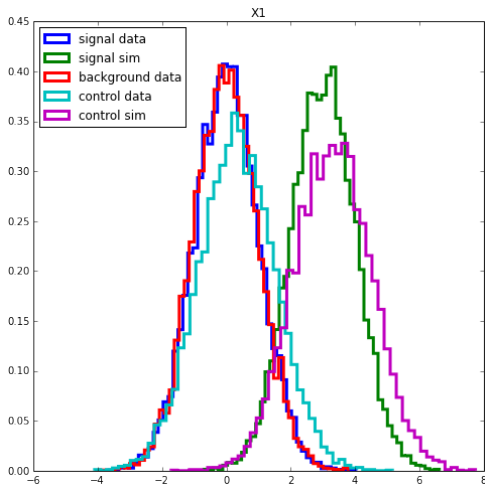
Therefore,

- The true performance of  $\varphi$  on real data may be significantly different (typically lower) than estimated on simulated signal events versus real data background events.
- Passing the KS test should not be interpreted as  $\varphi$  not exploiting simulation versus real data artefacts.

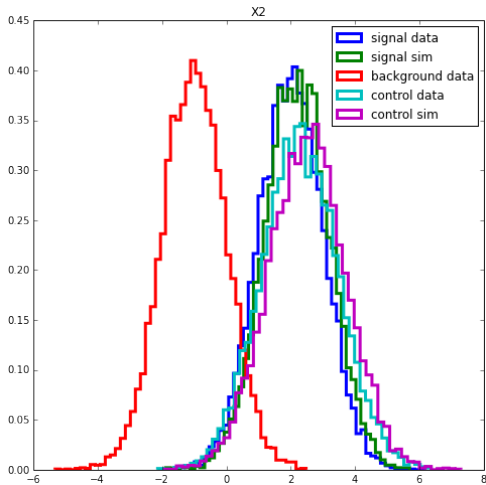
## Toy example

Let us consider an artificial classification problem between signal and background events, along with some close control channel data  $\mathcal{C}^{\text{sim}}$  and  $\mathcal{C}^{\text{data}}$ .

Let us assume an input space defined on three input variables  $X_1$ ,  $X_2$ ,  $X_3$  as follows.

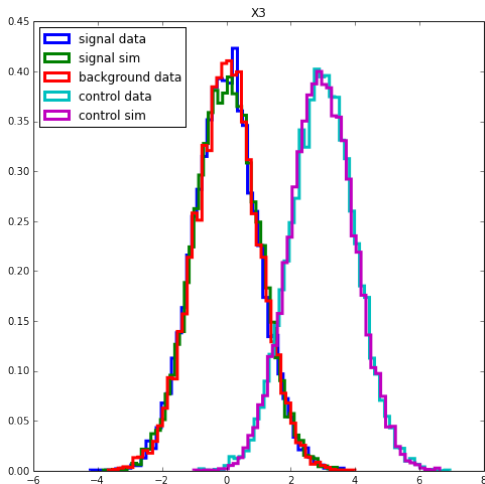


$X_1$  is **irrelevant** for real data signal versus real data background, but relevant for simulated versus real data events.



$X_2$  is **relevant** for background and non-background events.





$X_3$  is relevant for training versus control events, but has otherwise no discriminative power between signal and background events.

# Random exploration

```
def find_best_tree(X_train, y_train, X_test, y_test,
                  X_data, y_data, X_control_sim, X_control_data):
    best_auc_test, best_auc_data = 0, 0
    best_ks = 0
    best_tree = None

    for seed in range(2000):
        clf = ExtraTreesClassifier(n_estimators=1, max_features=1,
                                   max_leaf_nodes=5, random_state=seed)
        clf.fit(X_train, y_train)
        auc_test = roc_auc_score(y_test, clf.predict_proba(X_test)[: , 1])
        auc_data = roc_auc_score(y_data, clf.predict_proba(X_data)[: , 1])
        ks = ks_statistic(clf.predict_proba(X_control_sim)[: , 1],
                          clf.predict_proba(X_control_data)[: , 1])

        if auc_test > best_auc_test and ks < 0.09:
            best_auc_test = auc_test
            best_auc_data = auc_data
            best_ks = ks
            best_tree = clf

    return best_auc_test, best_auc_data, best_ks, best_tree
```

# Random exploration

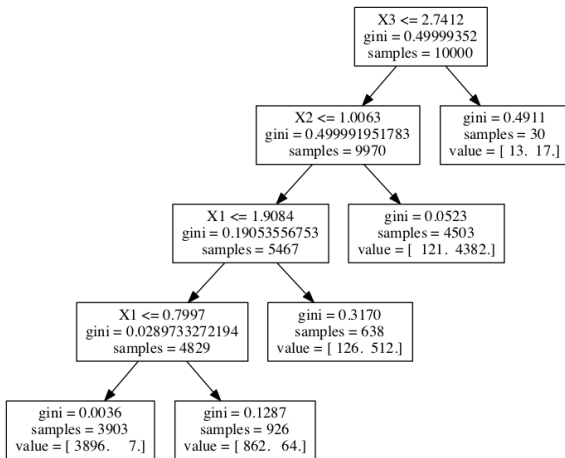
```
auc_test, auc_data, ks, tree = find_best_tree(...)
>>> auc_test = 0.9863 # Estimated AUC (simulated signal vs. data background)
>>> auc_data = 0.9097 # True AUC (data signal vs. data background)
>>> ks = 0.0578      # KS statistic < 0.09
```

What just happened? By chance, we have found a classifier that

- has seemingly good test performance;
- passes the control channel test that we have defined.

This classifier appears to be exactly the one we were seeking.

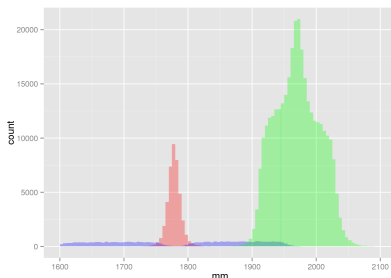
**Wrong.** The expected ROC AUC of 0.91 on real data signal and real data background is significantly lower than our first estimate, suggesting that there is still something wrong.



$\varphi$  exploits  $X_1$ , i.e. simulation versus real data artefacts to indirectly classify signal from background events, **while still passing the control channel test** because of its use of  $X_3$ !

# Winning the challenge

1. Learn to distinguish between training and control data,
2. Build a classifier on training data, with all the freedom to exploit simulation artefacts,
3. Assign random predictions to samples predicted as control data, otherwise predict using the classifier found in the previous step.



*The reconstructed mass allows to distinguish **signal** from **background** and training from **control**!*

## A machine learning response

As simulated training data increases (i.e., as  $N \rightarrow \infty$ ),

$$\frac{1}{N} \sum_{\mathbf{x}_i} L(\varphi(\mathbf{x}_i)) \rightarrow \int L(\varphi(\mathbf{x})) p_{\text{sim}}(\mathbf{x}) d\mathbf{x}.$$

We want to be good on real data, i.e., minimize

$$\int L(\varphi(\mathbf{x})) p_{\text{data}}(\mathbf{x}) d\mathbf{x}.$$

Solution: *importance weighting*.

$$\varphi^* = \arg \min_{\varphi} \frac{1}{N} \sum_{\mathbf{x}_i} \frac{p_{\text{data}}(\mathbf{x}_i)}{p_{\text{sim}}(\mathbf{x}_i)} L(\varphi(\mathbf{x}_i))$$

# Density ratio estimation

But for signal events, we don't even have real data observations!

$$\text{Assumption: } \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{sim}}(\mathbf{x})} \approx \frac{p_{\text{data-control}}(\mathbf{x})}{p_{\text{sim-control}}(\mathbf{x})} = r(\mathbf{x})$$

In the likelihood-free setting, estimating  $r(\mathbf{x})$  is known as the *density-ratio estimation* problem. Same as

- Learning under covariate shift,
- Probabilistic classification,
- Likelihood-ratio test,
- Outlier detection,
- Mutual information estimation, ...

See Sugiyama et al. (2012) for a review.

# Conclusions

- Formulating appropriate machine learning tasks is **difficult**.
- On purpose or **unwillingly**, simulation versus real data artefacts could be exploited to maximize classifiers accuracy.
- Physically more correct classifiers can be obtained e.g. with **density-ratio reweighting**.



# References

- Blake, T., Bettler, M.-O., Chrzaszcz, M., Dettori, F., Ustyuzhanin, A., and Likhomanenko, T. (2015). Flavours of physics: the machine learning challenge or the search of  $\tau \rightarrow 3\mu$  decays at LHCb.
- Louppe, G. and Head, T. (2015). Pitfalls of evaluating a classifier's performance in high energy physics applications.  
<http://dx.doi.org/10.5281/zenodo.34631>.
- Sugiyama, M., Suzuki, T., and Kanamori, T. (2012). *Density ratio estimation in machine learning*. Cambridge University Press.