

Approximating Likelihood Ratios with Calibrated Discriminative Classifiers

Kyle Cranmer and Gilles Louppe
New York University

January 22, 2016

Abstract

In particle physics, likelihood ratio tests are established tools for statistical inference. In many cases, these tests are complicated by the fact that computer simulators are used as a generative model for the data, which does not provide a way to evaluate the likelihood function. In this paper, we demonstrate that likelihood ratios are invariant under dimensionality reduction, provided the transformation is monotonic with the likelihood ratio. As a direct consequence, we show why and how discriminative classifiers can be used to approximate the likelihood ratio when only a generative model for the data is available. In particular, the proposed method offers a machine learning-based approach to statistical inference that is complementary to likelihood-free Bayesian inference algorithms, like Approximate Bayesian Computation, as it does not require the definition of a prior over model parameters. [GL: Also mention conclusions on experimental results, as required by JASA guidelines.]

Keywords: likelihood ratio, likelihood-free inference, classification, particle physics

1 Introduction

The likelihood function is the central object that summarizes the information from an experiment needed for inference of model parameters. The likelihood function is key to many areas of science that report the results of classical hypothesis tests or confidence intervals using the (generalized or profile) likelihood ratio as a test statistic. At the same time, with the advance of computing technology, it has become increasingly common that a simulator (or generative model) is used to describe complex processes that tie parameters θ of an underlying theory and measurement apparatus to high-dimensional observations \mathbf{x} . However, directly evaluating the likelihood function in these cases is often impossible or is computationally impractical. In this likelihood-free setting, various methods for statistical inference have been proposed, including most notably the Approximate Bayesian Computation (ABC) class of algorithms aiming at sampling from the posterior distribution of model parameters [\[GL: Add references\]](#).

The main result of this paper is to show that the likelihood ratio is invariant under dimensionality reduction, under the assumption that the corresponding transformation is monotonic with the likelihood ratio. As a direct consequence, we derive and propose an alternative machine learning-based approach for likelihood-free inference that can also be used in a classical (frequentist) setting where a prior over the model parameters is not available. More specifically, we demonstrate why and how discriminative classifiers can be used to construct equivalent likelihood ratio tests when only a generative model for the data is available for training and calibration.

As a concrete example, let us consider searches for new particles at the Large Hadron Collider (LHC). The simulator that is sampling from $p(\mathbf{x}|\theta)$ is based on quantum field theory, a detailed simulation of the particle detector, and data processing algorithms that

transform raw sensor data into the feature vector \mathbf{x} (Sjostrand et al., 2006; Agostinelli et al., 2003). The ATLAS and CMS experiments have published hundreds of papers where the final result was formulated as a hypothesis test or confidence interval using a generalized likelihood ratio test (Cowan et al., 2010), including most notably the discovery of the Higgs boson (The ATLAS Collaboration, 2012; The CMS Collaboration, 2012) and subsequent measurement of its properties. The bulk of the likelihood ratio tests at the LHC are based on the distribution of a single event-level feature that discriminates between a hypothesized process of interest (labeled *signal*) and various other processes (labeled *background*). Typically, pseudo-data from the simulator are used to approximate the density at various parameter points, and an interpolation algorithm is used to approximate the parameterized model (Cranmer et al., 2012). In particular, to improve the statistical power of these tests, hundreds of these searches have already been using supervised learning to train discriminative classifiers that take advantage of a high dimensional feature vector \mathbf{x} . [GL: We should then elaborate on how the proposed method could help these searches.]

The rest of this paper is organized as follows. In Section... [GL: todo]

2 Likelihood ratio tests

Let \mathbf{X} be a random vector with values $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^p$ and let $p_{\mathbf{X}}(\mathbf{x}|\theta)$ denote the density probability of \mathbf{X} at value \mathbf{x} under the parameterization θ . Let also assume some i.i.d. observed data $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$. In the setting where one is interested in simple hypothesis testing between a null $\theta = \theta_0$ against an alternate $\theta = \theta_1$, the Neyman-Pearson lemma states that the likelihood ratio

$$T(\mathcal{D}; \theta_0, \theta_1) = \prod_{\mathbf{x} \in \mathcal{D}} \frac{p_{\mathbf{X}}(\mathbf{x}|\theta_0)}{p_{\mathbf{X}}(\mathbf{x}|\theta_1)} \quad (2.1)$$

is the most powerful test statistic.

In order to evaluate $T(\mathcal{D})$, one must be able to evaluate the probability density $p_{\mathbf{X}}(\mathbf{x}|\theta)$ at any value \mathbf{x} . However, it is increasingly common in science that one has a complex simulation that can act as generative model for $p_{\mathbf{X}}(\mathbf{x}|\theta)$, but one cannot evaluate the density directly. For instance, this is the case in high energy physics (Neal, 2007) where the simulation of particle detectors can only be done in the forward mode. [GL: Reinsert citations for Scott and Nowak (2005) and Xin Tong (2013).]

3 Approximating likelihood ratios with classifiers

[GL: Think of a better section title? The equivalence is strict, not approximate.]

The main result of this paper (see Theorem 1) is to generalize the observation that one can form a test statistic

$$T'(\mathcal{D}; \theta_0, \theta_1) = \prod_{\mathbf{x} \in \mathcal{D}} \frac{p_{\mathbf{Y}}(y = s(\mathbf{x})|\theta_0)}{p_{\mathbf{Y}}(y = s(\mathbf{x})|\theta_1)} \quad (3.1)$$

that is strictly equivalent to 2.1, provided the change of variable $\mathbf{Y} = s(\mathbf{X})$ is based on a (parameterized) function s that is strictly monotonic with the density ratio

$$r(\mathbf{x}; \theta_0, \theta_1) = \frac{p_{\mathbf{X}}(\mathbf{x}|\theta_0)}{p_{\mathbf{X}}(\mathbf{x}|\theta_1)}. \quad (3.2)$$

As derived below, this allows to recast the original likelihood ratio test into an alternate form in which supervised learning can be used to build $s(\mathbf{x})$ as a discriminative classifier. In Section 4 we extend this result to generalized likelihood ratio tests, where it will be useful to have the classifier parameterized in terms of (θ_0, θ_1) .

3.1 Likelihood ratios under change of variables

Theorem 1. *Let \mathbf{X} be a random vector with values in $\mathcal{X} \subseteq \mathbb{R}^p$ and parameterized probability density $p_{\mathbf{X}}(\mathbf{x} = (x_1, \dots, x_p)|\theta)$ and let $s : \mathbb{R}^p \mapsto \mathbb{R}$ be a function monotonic with*

the density ratio $r(\mathbf{x}; \theta_0, \theta_1) = \frac{p_{\mathbf{X}}(\mathbf{x}|\theta_0)}{p_{\mathbf{X}}(\mathbf{x}|\theta_1)}$, for given parameters θ_0 and θ_1 . In these conditions,

$$r(\mathbf{x}; \theta_0, \theta_1) = \frac{p_{\mathbf{X}}(\mathbf{x}|\theta_0)}{p_{\mathbf{X}}(\mathbf{x}|\theta_1)} = \frac{p_{\mathbf{Y}}(y = s(\mathbf{x})|\theta_0)}{p_{\mathbf{Y}}(y = s(\mathbf{x})|\theta_1)}, \quad (3.3)$$

where $p_{\mathbf{Y}}(y = s(\mathbf{x}; \theta_0, \theta_1)|\theta)$ is the induced probability density of $\mathbf{Y} = s(\mathbf{X}; \theta_0, \theta_1)$.

Proof. Starting from the definition of the probability density function, we have

$$\begin{aligned} p_Y(y = s(\mathbf{x})|\theta_0) &= \frac{d}{dy} P(s(\mathbf{X}) \leq y) \\ &= \frac{d}{dy} \int_{\{\mathbf{x}' \in \mathbb{R}^p : s(\mathbf{x}') \leq y\}} p_{\mathbf{X}}(\mathbf{x}'|\theta_0) d\mathbf{x}' \\ &= \frac{d}{dy} \int_{\mathbb{R}^p} H(y - s(\mathbf{x}')) p_{\mathbf{X}}(\mathbf{x}'|\theta_0) d\mathbf{x}' \\ &= \int_{\mathbb{R}^p} \frac{d}{dy} H(y - s(\mathbf{x}')) p_{\mathbf{X}}(\mathbf{x}'|\theta_0) d\mathbf{x}' \\ &= \int_{\mathbb{R}^p} \delta(y - s(\mathbf{x}')) p_{\mathbf{X}}(\mathbf{x}'|\theta_0) d\mathbf{x}' \end{aligned} \quad (3.4)$$

where H and δ are respectively the Heaviside step and the Dirac delta functions. Intuitively, this last expression can be understood as the integral over all $\mathbf{x}' \in \mathbb{R}^p$ such that $s(\mathbf{x}') = y$, as picked by the Dirac delta function. Given Theorem 6.1.5 of Hörmander (1990), it further comes

$$\begin{aligned} p_Y(y = s(\mathbf{x})|\theta_0) &= \int_{\mathbb{R}^p} \delta(y - s(\mathbf{x}')) p_{\mathbf{X}}(\mathbf{x}'|\theta_0) d\mathbf{x}' \\ &= \int_{\{\mathbf{x}' \in \mathbb{R}^p : s(\mathbf{x}') = y\}} \frac{1}{|\nabla s(\mathbf{x}')|} p_{\mathbf{X}}(\mathbf{x}'|\theta_0) dS_{\mathbf{x}'} \end{aligned} \quad (3.5)$$

where $|\nabla s(\mathbf{x}')| = \sqrt{\sum_{i=1}^p \left| \frac{\partial}{\partial x_i} s(\mathbf{x}') \right|^2}$ and where $dS_{\mathbf{x}'}$ is the Euclidean surface measure on $\{\mathbf{x}' \in \mathbb{R}^p : s(\mathbf{x}') = y\}$. Also, since $s(\mathbf{x})$ is monotonic with $\frac{p_{\mathbf{X}}(\mathbf{x}|\theta_0)}{p_{\mathbf{X}}(\mathbf{x}|\theta_1)}$, it exists an invertible

function $m : \mathbb{R}^+ \mapsto \mathbb{R}$ such that $s(\mathbf{x}) = m\left(\frac{p_{\mathbf{X}}(\mathbf{x}|\theta_0)}{p_{\mathbf{X}}(\mathbf{x}|\theta_1)}\right)$. In particular, we have

$$\begin{aligned}\frac{p_{\mathbf{X}}(\mathbf{x}|\theta_0)}{p_{\mathbf{X}}(\mathbf{x}|\theta_1)} &= m^{-1}(s(\mathbf{x})) \\ p_{\mathbf{X}}(\mathbf{x}|\theta_0) &= m^{-1}(s(\mathbf{x}))p_{\mathbf{X}}(\mathbf{x}|\theta_1)\end{aligned}\tag{3.6}$$

Combining equations 3.5 and 3.6, we have

$$\begin{aligned}p_Y(y = s(\mathbf{x})|\theta_0) &= \int_{\{\mathbf{x}' \in \mathbb{R}^p : s(\mathbf{x}') = y\}} \frac{p_{\mathbf{X}}(\mathbf{x}'|\theta_0)}{|\nabla s(\mathbf{x}')|} dS_{\mathbf{x}'} \\ &= \int_{\{\mathbf{x}' \in \mathbb{R}^p : s(\mathbf{x}') = y\}} \frac{1}{|\nabla s(\mathbf{x}')|} m^{-1}(s(\mathbf{x}')) p_{\mathbf{X}}(\mathbf{x}'|\theta_1) dS_{\mathbf{x}'} \\ &= \int_{\{\mathbf{x}' \in \mathbb{R}^p : s(\mathbf{x}') = y\}} \frac{1}{|\nabla s(\mathbf{x}')|} m^{-1}(y) p_{\mathbf{X}}(\mathbf{x}'|\theta_1) dS_{\mathbf{x}'} \\ &= m^{-1}(s(\mathbf{x})) \int_{\{\mathbf{x}' \in \mathbb{R}^p : s(\mathbf{x}') = y\}} \frac{1}{|\nabla s(\mathbf{x}')|} p_{\mathbf{X}}(\mathbf{x}'|\theta_1) dS_{\mathbf{x}'} \\ &= \frac{p_{\mathbf{X}}(\mathbf{x}|\theta_0)}{p_{\mathbf{X}}(\mathbf{x}|\theta_1)} \int_{\{\mathbf{x}' \in \mathbb{R}^p : s(\mathbf{x}') = y\}} \frac{1}{|\nabla s(\mathbf{x}')|} p_{\mathbf{X}}(\mathbf{x}'|\theta_1) dS_{\mathbf{x}'}\end{aligned}\tag{3.7}$$

Similarly, Equation 3.5 can be used to derive $p_Y(y = s(\mathbf{x})|\theta_1)$, finally yielding

$$\begin{aligned}\frac{p_Y(y = s(\mathbf{x})|\theta_0)}{p_Y(y = s(\mathbf{x})|\theta_1)} &= \frac{p_{\mathbf{X}}(\mathbf{x}|\theta_0)}{p_{\mathbf{X}}(\mathbf{x}|\theta_1)} \frac{\int_{\{\mathbf{x}' \in \mathbb{R}^p : s(\mathbf{x}') = y\}} \frac{1}{|\nabla s(\mathbf{x}')|} p_{\mathbf{X}}(\mathbf{x}'|\theta_1) dS_{\mathbf{x}'}}{\int_{\{\mathbf{x}' \in \mathbb{R}^p : s(\mathbf{x}') = y\}} \frac{1}{|\nabla s(\mathbf{x}')|} p_{\mathbf{X}}(\mathbf{x}'|\theta_1) dS_{\mathbf{x}'}} \\ &= \frac{p_{\mathbf{X}}(\mathbf{x}|\theta_0)}{p_{\mathbf{X}}(\mathbf{x}|\theta_1)}.\end{aligned}\tag{3.8}$$

□

3.2 Probabilistic classification for likelihood ratios

Proposition 2. *Let $\mathbf{X} = (X_1, \dots, X_p)$ and Y be random input and output variables with values in $\mathcal{X} \subseteq \mathbb{R}^p$ and $\mathcal{Y} = \{0, 1\}$ and mixed joint probability density function $p_{\mathbf{X}, Y}(\mathbf{x}, y)$.*

For the squared error loss, the best regression function $s : \mathcal{X} \mapsto [0, 1]$, or equivalently the best probabilistic classifier, is

$$s^*(\mathbf{x}) = \frac{P(Y = 1)p_{\mathbf{X}|Y}(\mathbf{x}|Y = 1)}{P(Y = 0)p_{\mathbf{X}|Y}(\mathbf{x}|Y = 0) + P(Y = 1)p_{\mathbf{X}|Y}(\mathbf{x}|Y = 1)}. \quad (3.9)$$

Proof. For the squared error loss,

$$\begin{aligned} s^*(\mathbf{x}) &= \arg \min_{s(\mathbf{x})} \mathbb{E}_{Y|\mathbf{X}=\mathbf{x}} \{(Y - s(\mathbf{x}))^2\} \\ &= \arg \min_{s(\mathbf{x})} \mathbb{E}_{Y|\mathbf{X}=\mathbf{x}} \{Y^2\} - 2s(\mathbf{x})\mathbb{E}_{Y|\mathbf{X}=\mathbf{x}} \{Y\} + s(\mathbf{x})^2 \\ &= \arg \min_{s(\mathbf{x})} -2s(\mathbf{x})\mathbb{E}_{Y|\mathbf{X}=\mathbf{x}} \{Y\} + s(\mathbf{x})^2 \end{aligned} \quad (3.10)$$

The last expression is minimized when $\frac{d}{ds(\mathbf{x})}(-2s(\mathbf{x})\mathbb{E}_{Y|\mathbf{X}=\mathbf{x}} \{Y\} + s(\mathbf{x})^2) = 0$, that is when $-2\mathbb{E}_{Y|\mathbf{X}=\mathbf{x}} \{Y\} + 2s(\mathbf{x}) = 0$, hence

$$s^*(\mathbf{x}) = \mathbb{E}_{Y|\mathbf{X}=\mathbf{x}} \{Y\}. \quad (3.11)$$

For $\mathcal{Y} = \{0, 1\}$,

$$\begin{aligned} \mathbb{E}_{Y|\mathbf{X}=\mathbf{x}} \{Y\} &= P(Y = 0|\mathbf{X} = \mathbf{x}) \times 0 + P(Y = 1|\mathbf{X} = \mathbf{x}) \times 1 \\ &= \frac{P(Y = 1)p_{\mathbf{X}|Y}(\mathbf{x}|Y = 1)}{p_{\mathbf{X}}(\mathbf{x})} \\ &= \frac{P(Y = 1)p_{\mathbf{X}|Y}(\mathbf{x}|Y = 1)}{P(Y = 0)p_{\mathbf{X}|Y}(\mathbf{x}|Y = 0) + P(Y = 1)p_{\mathbf{X}|Y}(\mathbf{x}|Y = 1)}. \end{aligned} \quad (3.12)$$

□

For $P(Y = 0) = P(Y = 1) = \frac{1}{2}$, the best regression function s^* simplifies to

$$s^*(\mathbf{x}) = \frac{p_{\mathbf{X}|Y}(\mathbf{x}|Y = 1)}{p_{\mathbf{X}|Y}(\mathbf{x}|Y = 0) + p_{\mathbf{X}|Y}(\mathbf{x}|Y = 1)}. \quad (3.13)$$

If we further assume that samples for $Y = 0$ are drawn from some parameterized distribution with probability density $p_{\mathbf{X}}(\mathbf{x}|\theta_0)$ and similarly for $Y = 1$, then the best regression function can be rewritten as

$$s^*(\mathbf{x}) = \frac{p_{\mathbf{X}}(\mathbf{x}|\theta_1)}{p_{\mathbf{X}}(\mathbf{x}|\theta_0) + p_{\mathbf{X}}(\mathbf{x}|\theta_1)}. \quad (3.14)$$

In particular, this regression function satisfies conditions of Theorem 1 since $s^*(\mathbf{x}) = m(\frac{p_{\mathbf{X}}(\mathbf{x}|\theta_0)}{p_{\mathbf{X}}(\mathbf{x}|\theta_1)})$, for $m(r(\mathbf{x})) = \frac{1}{1+r(\mathbf{x})}$, is monotonic with $\frac{p_{\mathbf{X}}(\mathbf{x}|\theta_0)}{p_{\mathbf{X}}(\mathbf{x}|\theta_1)}$. In other words, Proposition 2 yields a sufficient procedure for Theorem 1 to hold, guaranteeing that any *universally strongly consistent* algorithm can be used for learning s^* . Note however, that it is not a necessary procedure since Theorem 1 holds for any monotonic function m of the density ratio, i.e., not for $m(r(\mathbf{x})) = \frac{r(\mathbf{x})}{1+r(\mathbf{x})}$ only.

Equivalently, since $m^{-1}(s^*(\mathbf{x})) = \frac{1-s^*(\mathbf{x})}{s^*(\mathbf{x})} = r(\mathbf{x})$, Proposition 2 shows that the likelihood ratio estimation problem is directly related to probabilistic classification, and that algorithms proposed to solve one can be used for solving the other. In this context, Theorem 1 shows that in case we learn a probabilistic classifier $s(\mathbf{x})$ which is imperfect up to a monotonic transformation of $r(\mathbf{x})$, then one can still resort to calibration (i.e., modeling $p_Y(y = s(\mathbf{x}))$) to compute $r(\mathbf{x})$ exactly.

3.3 Approximating the reduction map and its distribution

In order for this approach to be useful in the likelihood-free setting, we need to be able to approximate both $s(\mathbf{x})$ and $p(s(\mathbf{x})|\theta)$ based on a finite number of samples $\{\mathbf{x}_i\}$ drawn from the generative model $p(\mathbf{x}|\theta)$. Let denote the approximated dimensionality reduction map $\hat{s}(\mathbf{x})$ and its approximated distribution $\hat{p}(\hat{s}(\mathbf{x})|\theta)$. [GL: Explain how to approximate $p(s)$ in practice?]

One strength of this approach is that it factorizes the approximation of the per-event

likelihood ratio ($\hat{s}(\mathbf{x}) \approx s(\mathbf{x})$) from the calibration procedure ($\hat{p}(\hat{s}(\mathbf{x})|\theta) \approx p(\hat{s}(\mathbf{x})|\theta)$). Thus, even if the classifier does a poor job at reproducing the level sets of the per-event likelihood ratio, the density of \hat{s} can still be well calibrated. In that case, one might loose power, but the resulting inference will still be valid. [GL: Shouldn't we elaborate?] This point was made by Neal (2007) and is well appreciated by the particle physics community that typically takes a conservative attitude towards the use of machine learning classifiers precisely due to concerns about the calibration p -values in the face of nuisance parameters associated to the simulator.

[GL: We should comment as to what happens if $\hat{s}(x)$ is not monotonic with r]

[GL: Add pseudo-code to summarize the whole pipeline?]

3.4 Classification and frequentist hypothesis tests

[GL: Shall we keep everything from this?]

Vast literature exists around generative and discriminative classifiers (Ng and Jordan, 2002). Typically, generative classifiers learn a model for the joint probability $p(\mathbf{x}, y)$, of the inputs x and the classification label y , and predict $p(y|\mathbf{x})$ via Bayes rule. In contrast, discriminative classifiers model the posterior $p(y|\mathbf{x})$ directly. For classification tasks, one then thresholds on $p(y|\mathbf{x})$. In both cases this description in terms of a posterior requires a prior distribution for $p(y)$, which is either modeled explicitly or learned from the training data. This familiar formulation of classification may lead to some confusion in the setting of the current work.

The first possible source of confusion we wish to avoid is that here $p(\mathbf{x}|\theta)$ is a *generative statistical model* for the features \mathbf{x} , not a generative classifier. We think of the $p(x|\theta)$ along the lines of a traditional scientific theory or simulator, able to make predictions about \mathbf{x} and being motivated by domain-specific considerations.

The second possible source of confusion is that we are not directly interested in calibrating the classification score in terms of a per-event posterior probability $p(y|\mathbf{x})$. Instead, we are interested in the approximation of the per-experiment likelihood function or likelihood ratio, which might be used for several purposes, including the calculation of p-values.

Lastly, in the setting of frequentist hypothesis tests and confidence intervals, we do not have a prior $\pi(\theta)$. While we can use the generative models to produce training data $(\mathbf{x}, y = 0)$ generated from $p(\mathbf{x}|\theta_0)$ and $(\mathbf{x}, y = 1)$ generated from $p(\mathbf{x}|\theta_1)$, the relative mix $p(y)$ is arbitrary. Since the prior $p(y)$ is not needed for the target likelihood ratio test and because the classifier score $p(y|\mathbf{x})$ may not be well calibrated, we choose to denote the classifier score $s(\mathbf{x})$ and simply think of it as a deterministic dimensionality reduction map $s : \mathbb{R}^p \rightarrow \mathbb{R}$. Similar points have been made by Scott and Nowak (2005) and Neal (2007).

4 Generalized likelihood ratio tests

Thus far we have shown that the target likelihood ratio $r(\mathbf{x}; \theta_0, \theta_1) = p(\mathbf{x}|\theta_0)/p(\mathbf{x}|\theta_1)$ with high dimensional features \mathbf{x} can be reproduced via the univariate densities $p(s(\mathbf{x})|\theta_0)$ and $p(s(\mathbf{x})|\theta_1)$ if the reduction $s(\mathbf{x})$ is monotonic with $r(\mathbf{x}; \theta_0, \theta_1)$. We now generalize from the ratio of two simple hypotheses specified by θ_0 and θ_1 to the case of composite hypothesis testing where θ are continuous model parameters.

In the case of composite hypotheses $\theta \in \Theta_0$ against an alternative $\theta \in \Theta_1$ (possibly such that $\Theta_0 \subseteq \Theta_1$), the generalized likelihood ratio test, also known as the profile likelihood ratio test, is commonly used

$$\Lambda(\mathcal{D}; \Theta_0; \Theta_1) = \frac{\sup_{\theta \in \Theta_0} p(\mathcal{D}|\theta)}{\sup_{\theta \in \Theta_1} p(\mathcal{D}|\theta)} . \quad (4.1)$$

This generalized likelihood ratio can be used both for hypothesis tests in the presence of nuisance parameters or to create confidence intervals with or without nuisance param-

ters. Often, the parameter vector is broken into two components $\theta = (\mu, \nu)$, where the μ components are considered parameters of interest while the ν components are considered nuisance parameters. In that case Θ_0 corresponds to all values of ν with μ fixed.

Evaluating the generalized likelihood ratio as defined by Eqn. 4.1 requires finding for both the numerator and the denominator the maximum likelihood estimator

$$\hat{\theta} = \arg \max_{\theta} p(\mathcal{D}|\theta). \quad (4.2)$$

Again, this is made difficult in the likelihood-free setting and it is not obvious that we can find the same estimators if we are working instead with $p(s(\mathbf{x})|\theta)$. Fortunately, there is a construction based on s that works: the maximum likelihood estimate of Eqn. 4.2 is the same as the value that maximizes the likelihood ratio with respect to $p(\mathcal{D}|\theta_1)$, for some fixed value of θ_1 . This allows us to use Theorem 1 to reformulate the maximum likelihood estimate as

$$\begin{aligned} \hat{\theta} &= \arg \max_{\theta} p(\mathcal{D}|\theta) \\ &= \arg \max_{\theta} \prod_{\mathbf{x} \in \mathcal{D}} \frac{p(\mathbf{x}|\theta)}{p(\mathbf{x}|\theta_1)} \\ &= \arg \max_{\theta} \prod_{\mathbf{x} \in \mathcal{D}} \frac{p(s(\mathbf{x}; \theta, \theta_1)|\theta)}{p(s(\mathbf{x}; \theta, \theta_1)|\theta_1)}, \end{aligned} \quad (4.3)$$

where $s(\mathbf{x}; \theta, \theta_1)$ denotes a *parameterized* transformation s of \mathbf{X} in terms of (θ, θ_1) that is monotonic with $r(\mathbf{x}; \theta, \theta_1)$. Note that it is important that we include the denominator $p(s(\mathbf{x}; \theta, \theta_1)|\theta_1)$ because this cancels Jacobian factors that vary with θ .

4.1 Parameterized classifier

In order to provide parameter inference in the likelihood-free setting as described above, we must train a family $s(\mathbf{x}; \theta, \theta_1)$ of classifiers parameterized by θ_0 and θ_1 , the parameters

Algorithm 1 Training of a parameterized classifier.

```
 $\mathcal{L} := \{\}$   
for  $\theta_0$  in  $\Theta_0$  do  
  for  $\theta_1$  in  $\Theta_1$  do  
    generate  $\mathbf{x} \sim p(\mathbf{x}|\theta_0)$   
    append  $\{(\mathbf{x}, \theta_0, \theta_1, y = 0)\}$  to  $\mathcal{L}$   
    generate  $\mathbf{x} \sim p(\mathbf{x}|\theta_1)$   
    append  $\{(\mathbf{x}, \theta_0, \theta_1, y = 1)\}$  to  $\mathcal{L}$   
  end for  
end for  
Learn  $s(\mathbf{x}; \theta_0, \theta_1)$  from  $\mathcal{L}$ 
```

associated to the null and alternate hypotheses, respectively. While this could be done independently for all θ_0 and θ_1 , using the procedure outlined in Section 3.2, it is desirable and convenient to have a smooth evolution of the classification score as a function of the parameters. For this reason, we anticipate a single learning stage based on training data with input $(\mathbf{x}, \theta_0, \theta_1)_i$ and target y_i , as outlined in Algorithm 1. Somewhat unusually, the unknown values of the parameters are taken as input to the classifier; their values will be specified via the enveloping (generalized) likelihood ratio of Eqn. 4.1.

While the function $p(\mathbf{x}|\theta_1)/(p(\mathbf{x}|\theta_0) + p(\mathbf{x}|\theta_1))$ will minimize the expected squared error loss based on training data produced according to Algorithm 1, it is not clear how training data from $\theta'_0 \neq \theta_0$ and $\theta'_1 \neq \theta_1$ will influence a real world classifier with finite capacity. This is left as an area for future work.

[GL: Shall we keep Algorithm 1 or rather give a sketch in words? In its current state, the algorithm does not explain how to build \mathcal{L} when the parameter spaces are continuous.

It also does not specify how many samples should be generated per parameter value.]

4.2 Parameterized calibration

Once the classifier is trained, we can use the generative model together with a univariate density estimation technique (e.g. histograms or kernel density estimation) to approximate $p(s|\theta)$ for specific parameter points. For a single parameter point, this is a tractable univariate density estimation problem. The challenge comes from the need to calibrate this density for all values of θ . A straight forward approach would be to run the generative model on demand for any particular value of θ . In the context of a likelihood fit this would mean that the optimization algorithm that is trying to maximize the likelihood with respect to θ needs access to the generative model $p(\mathbf{x}|\theta)$. This can be impractical when the generative model is computationally expensive or has high-latency (for instance some human intervention is required to reconfigure the generative model). In HEP, with a fixed classifier, it has become common to interpolate the distribution between discrete values of θ in order to produce a continuous parameterization for $p(s|\theta)$ (Cranmer et al., 2012). One can easily imagine a number of approaches to embedding the classifier and estimating the density $p(s|\theta)$ and the relative merits of those approaches will depend critically on the dimensionality of θ and the computational cost of the generative model. We leave a more general strategy for this overarching optimization problem as an area of future work.

[GL: Sub-sections 4.1 and 4.2 could directly be inlined within section 4. What do you think?]

5 Applications

5.1 High energy physics

[GL: Add a concrete HEP example where we illustrate and compare the method for several classifiers and calibration algorithms.] [GL: Notations need to be adapted.] [GL: Rethink the outline? Start with comparison of mixtures (5.1.4) and then specialize to background vs background+signal, and thereby discuss that we can focus the capacity of the classifier.]

In high energy physics, we are often searching for some class of events, generically referred to as *signal*, in the presence of a separate class of *background* events. For each event we measure some quantities x that have corresponding distributions $p_b(x|\nu)$ for background and $p_s(x|\nu)$ for signal, where ν are nuisance parameters describing uncertainties in the underlying physics prediction or response of the measurement device. The total model is a mixture of the signal and background, and μ is the mixture coefficient associate to the signal component.

$$p(D | \mu, \nu) = \prod_{e=1}^n [\mu p_s(x_e | \nu) + (1 - \mu) p_b(x_e | \nu)] , \quad (5.1)$$

New particle searches at the LHC are typically framed as hypothesis test where the null corresponds to $\mu = 0$, and the generalized likelihood ratio is used as a test statistic.

5.1.1 Typical usage pattern

In this setting, large samples of pseudo-data $\{x_i, y_i\}$ generated with some nominal values of the parameters ν_0 , where $y = 0$ corresponds to the background density $p_b(x|\nu_0)$ and $y = 1$ corresponds to signal density $p_s(x|\nu_0)$. Importantly, the $y = 1$ label corresponds to the signal only, and not to the alternate signal-plus-background hypothesis. The resulting classifier approximates the regression function $p_s(x|\nu_0)/(p_s(x|\nu_0) + p_b(x|\nu_0))$, which is one

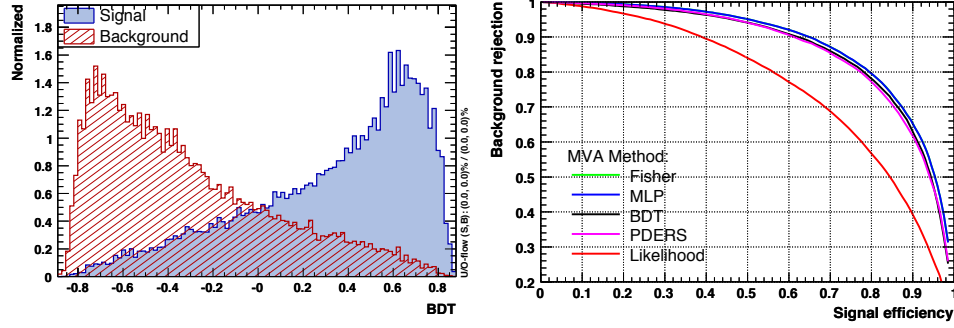


Figure 1: Left: an example of the distributions $p_b(\hat{s}|\nu)$ and $p_s(\hat{s}|\nu)$ when the classifier s is a boosted-decision tree (BDT). Right: the corresponding ROC curve (right) for this and other classifiers. (Figures taken from TMVA manual.)

to one with the likelihood ratio of the null to the alternate $p(x|\mu = 0, \nu_0)/p(x|\mu, \nu_0)$ for all μ . Associating the $y = 1$ label to the signal component has advantages because it helps the classifier focus its capacity on the relevant regions in the feature space, particularly when the signal is a very small perturbation to the background (i.e. $\mu \ll 1$).

Once the classifier is trained, large samples of pseudo-data are drawn from $p_s(x|\nu)$ and $p_b(x|\nu)$ and we estimate the distributions $\hat{p}_s(\hat{s}|\nu)$ and $\hat{p}_b(\hat{s}|\nu)$ continuously parameterized in ν . An example of the distributions of \hat{s} for the signal and background events with $\nu = \nu_0$ is shown in Figure 1.

These steps feed into a subsequent statistical test based on the observed data $D = (x_1, \dots, x_n)$. For each event, the classifier is evaluated and one performs inference on a parameter μ related to the presence of the signal contribution. In particular, one forms the statistical

model¹

$$p(D | \mu, \nu) = \prod_{e=1}^n [\mu \hat{p}_s(\hat{s}(x_e) | \nu) + (1 - \mu) \hat{p}_b(\hat{s}(x_e) | \nu)] . \quad (5.2)$$

5.1.2 Comments on typical usage of machine learning in HEP

Nuisance parameters are an after thought in the typical usage of machine learning in HEP. In fact, most discussions would related to the training and optimizing the classifier only consider $p_b(x)$ and $p_s(x)$ with $\nu = \nu_0$ being implicit. However, as experimentalists we know that we must account for various forms of systematic uncertainty, parameterized by nuisance parameters ν . In practice, we take the classifier as fixed and then propagate uncertainty through the classifier as in Eq. 5.2. Building the distribution $p(\hat{s}|\nu)$ for values of ν other than the nominal ν_0 used to train the classifier can be thought of as a calibration necessary for classical statistical inference; however, this classifier is clearly not optimal for $\nu \neq \nu_0$.

5.1.3 A more powerful approach

The standard use of machine learning in HEP can be improved by training a parameterized classifier corresponding to the generalized likelihood ratio test

$$\lambda(\mu) = \frac{p(D|\mu, \hat{\nu})}{p(D|\hat{\mu}, \hat{\nu})} , \quad (5.3)$$

following the approach outlined in Section ??.

There is an interesting distinction between this approach and the standard use in which the classifier is trained for a fixed ν_0 . In the standard use one trains a classifier for signal vs. background, which is equivalent (in an ideal setting) to training a classifier for

¹Sometimes there is an additional Poisson term when expected number of signal and background events is known, which is referred to as an extended likelihood or marked Poisson model.

null (background-only) vs. alternate (signal-plus-background) since the resulting regression functions are one-to-one with each other. In contrast, in the case of the generalized likelihood ratio test

$$\frac{p(x|0, \hat{\nu})}{p(x|\hat{\mu}, \hat{\nu})} = \frac{p_b(x|\hat{\nu})}{\hat{\mu}p_s(x_e|\hat{\nu}) + (1 - \hat{\mu})p_b(x_e|\hat{\nu})} , \quad (5.4)$$

the background components don't cancel and there is an additional term $p_b(x|\hat{\nu})/p_b(x|\hat{\nu})$. In practice, with classifiers of finite capacity, there will be some trade-off between taking into account this additional term and the more challenging learning problem when μ is very small.

5.1.4 Decomposing tests between mixture models into their components

In this section we generalize the capacity focusing technique of training classifiers to discriminate between components of a mixture model. First, we generalize Eq. 5.1 to a mixture model of several components

$$p(x|\theta) = \sum_c w_c(\theta) p_c(x|\theta) . \quad (5.5)$$

It is possible to re-write the target likelihood ratio between two mixture models in terms of pairwise classification problems.

$$\begin{aligned} \frac{p(x|\theta_0)}{p(x|\theta_1)} &= \frac{\sum_c w_c(\theta_0) p_c(x|\theta_0)}{\sum_{c'} w_{c'}(\theta_1) p_{c'}(x|\theta_1)} \\ &= \sum_c \left[\sum_{c'} \frac{w_{c'}(\theta_1)}{w_c(\theta_0)} \frac{p_{c'}(x|\theta_1)}{p_c(x|\theta_0)} \right]^{-1} \\ &= \sum_c \left[\sum_{c'} \frac{w_{c'}(\theta_1)}{w_c(\theta_0)} \frac{p_{c'}(s_{c,c',\theta_0,\theta_1}|\theta_1)}{p_c(s_{c,c',\theta_0,\theta_1}|\theta_0)} \right]^{-1} \end{aligned} \quad (5.6)$$

The second line is a trivial, but useful decomposition into pair-wise classification between $p_{c'}(x|\theta_1)$ and $p_c(x|\theta_0)$. The third line uses Theorem 1 to relate the high-dimensional likelihood ratio into an equivalent calibrated likelihood ratio based on the univariate density of the corresponding classifier, denoted $s_{c,c',\theta_0,\theta_1}$. In the situation where the only free parameters of the model are the mixture coefficients w_c , then the distributions $p_c(s_{c,c',\theta_0,\theta_1}|\theta)$ are independent of θ and can be pre-computed (after training the discriminative classifier, but before evaluating the likelihood ratio). Equation 5.6 allows one to take advantage of both the parameterized classifier as in Eq. 5.4 and the capacity focusing technique in the typical HEP usage pattern.

5.2 Likelihood-free inference

[GL: Rework this section, add a concrete example and compare with other likelihood-free inference algorithms.]

While the original motivation for this work was to improve the treatment of systematic uncertainties in new particle searches by parameterizing the classifier in terms of the nuisance parameters ν , the same approach can be used for parameter inference. In the case of new particle searches the parameter of interest is the mixture coefficient for the signal component $p_s(x|\nu)$. When measuring particle properties the distribution of the features also depend on parameters such as a particle's mass and quantum numbers. This is easily accommodated by extending $p_s(x|\nu) \rightarrow p_s(x|\theta)$, where θ includes both parameters of interest and nuisance parameters.

This formalism represents a significant step forward in the usage of machine learning in HEP, where classifiers have always been used between two static classes of events and not parameterized explicitly in terms of the physical quantities we wish to measure. The work of (Whiteson and Whiteson, 2007) is similar as the stochastic optimization was

directly trying to minimize the measurement uncertainty of a particle’s mass; however, the resulting classifier was fixed. This approach also offers the advantage that it explicitly reformulates the per-experiment optimization to the per-event optimization, which is less computationally intensive.

Another approach that is similar in spirit is the so-called matrix element method, in which one directly computes an approximate likelihood ratio by performing a computationally intensive integral associated to the detector response (Volobouev, 2011). In the approach considered in this paper, the detector response is naturally handled by the Monte Carlo sampling used in the simulation of the detector; however, that integral is intractable for the matrix element method. Even with drastic simplifications of the detector response, the matrix element method can take several minutes of CPU time to calculate the likelihood ratio $q(x)$ for a single event. The work here can be seen as aiming at the same conceptual target, but utilizing machine learning to overcome the complexity of the detector simulation. It also offers enormous speed increase for evaluating the likelihood at the cost of an initial training stage. In practice, the matrix element method has only been used for searches and measurement of a single physical parameter (sometimes with a single nuisance parameter as in (Aaltonen et al., 2010)).

Contemporary examples where the technique presented here could have major impact on HEP include the measurement of coefficients to quantum mechanical operators describing the decay of the Higgs boson (Chen et al., 2015) and, if we are so lucky, measurement of the mass of supersymmetric particles in cascade decays (Allanach et al., 2000). Both of these examples involve data sets with many events, each with a feature vector x that has on the order of 10 components, and a parameter vector θ with 5-10 parameters of interest and possibly many more nuisance parameters. The state of the art for the operator coefficients of the Higgs decay uses the so-called matrix element likelihood analysis (MELA) in which

the equivalent of $s(x; \theta_0, \theta_1)$ is approximated by neglecting detector effects (Gao et al., 2010; Bolognesi et al., 2012).

6 Related works

[GL: Add related works on density ratio estimation.] [GL: Remove others that are less directly related.]

Scott and Nowak (2005) and Xin Tong (2013) consider the machine learning problem associated to Neyman-Pearson hypothesis testing. As in this work, they consider the situation where one does not have access to the underlying distributions, but only has i.i.d. samples from each hypothesis. This work generalizes that goal from the Neyman-Pearson setting to generalized likelihood ratio tests and emphasizes the connection with classification. Perhaps a formal treatment similar to the Neyman-Pearson case can be brought to bear in this more general setting. In a similarly titled work, Gutmann et al. (2014) advocate using the cross-validated classification accuracy as the similarity metric used in ABC. While the goal there is also parameter inference in the likelihood-free setting, “classifier ABC” is very different than the approach presented here. Jaakkola and Haussler (1998) explore a way of leveraging generative models to derive kernel functions for use in discriminative methods. This interesting work is distinct from the point made here in which the generative model is being used for the purpose of providing training data and calibration. Zadrozny and Elkan (2001) emphasize the importance of calibrated probability estimates from decision trees and naive Bayesian classifiers and investigate various approaches to achieve this. In contrast to that work, we are not interested in calibrated probability estimates for $p(y|x)$ for individual events, but instead we use the calibration to correct for non-linear transformations of the target likelihood ratio and, perhaps, to provide calibrated

p-values based on those likelihood ratio tests. Ihler et al. (2004) take on a different problem (tests of statistical independence) by using machine learning algorithms to find scalar maps from the high-dimensional feature space that achieve the desired statistical goal when the fundamental high-dimensional test is intractable.

Neal (2007) also considered the problem of approximating the likelihood function when only a generative model is available. That work sketches a scheme in which one uses a classifier with both x and θ as an input to serve as a dimensionality reduction map. The key distinction comes in the handling of θ . Neal says “we cannot use the classifier on real data, since we don’t know the correct value for $[\theta]$ ” and goes on to outline an approach where one uses regression on a per-event basis to estimate $\hat{\theta}(x)$ and perform the composition $s(x; \hat{\theta}(x))$ (much like profiling).² This can lead to a significant loss of information since (at least in most particle physics examples) a single event carries little information about the true value of θ , though the full data set D may be informative – for instance, a single observation would not be sufficient to estimate the variance of a distribution, though repeated observations would. The work of Neal correctly identifies this as an approximation of the target likelihood even in the case of a ideal classifier. In contrast, the approach described here does not eliminate the dependence of the classifier on θ .³ Instead, we embed a parameterized classifier into the likelihood and postpone the evaluation of the classifier to the point of evaluation of the likelihood when θ is explicitly

²Neal considers a lower-dimensional ‘bottlenecks’ $\theta^* = g(\theta)$, which are not essential to the discussion here.

³As a technical point, in Neal’s work, the focus is on approximating the likelihood function (up to a multiplicative constant), which is equivalent to evaluating the ratio with respect to a fixed θ_1 as in Eq. 4.3. In Neal’s case, the dependence on θ is eliminated via $s(x; \hat{\theta}(x))$ and the map is constant; however, in this approach the map ratio is explicitly parameterized in terms of θ , so the ratio is important for canceling the corresponding Jacobian factors.

being tested. This avoids the loss of information that occurs from the regression step $\hat{\theta}(x)$ proposed by Neal and leads to Theorem 1, which is an exact result in the case of an ideal classifier. In both cases, the quality of the classifier is factorized from the calibration of its density, which allows for valid inference even if there is a loss of power due to a non ideal classifier.

7 Conclusions

We have outlined a technique to reformulate generalized likelihood ratio test over a high-dimensional data set with multiple events in terms of a univariate density of a classifier score. We have shown that a parameterized family of discriminative classifiers $\hat{s}(x; \theta_0, \theta_1)$ trained and calibrated with a simulator $p(x|\theta)$ can be used to approximate the likelihood ratio $\prod_i p(x_i|\theta_0)/p(x_i|\theta_1)$ even when it is not possible to directly evaluate the likelihood $p(x|\theta)$. It offers an alternative to approximate Bayesian computation for parameter inference in the likelihood-free setting that can also be used in the frequentist formalism without specifying a prior over the parameters. A strength of this approach is that it separates the quality of the approximation of the target likelihood from the quality of the calibration. The former is related to the ability of supervised learning approaches to classification, which will continue to improve. The calibration procedure for a particular parameter point is fairly straight forward since it involves estimating a univariate density using a generative model of the data. The difficulty of the calibration stage is performing this calibration continuously in θ . Different strategies to this calibration are anticipated depending on the dimensionality of θ , the complexity of the resulting likelihood function, and the difficulties associated to running the simulator.

References

- Aaltonen, T. et al. (2010). Top Quark Mass Measurement in the Lepton + Jets Channel Using a Matrix Element Method and *in situ* Jet Energy Calibration. *Phys.Rev.Lett.*, 105:252001.
- Agostinelli, S. et al. (2003). GEANT4: A Simulation toolkit. *Nucl.Instrum.Meth.*, A506:250–303.
- Allanach, B., Lester, C., Parker, M. A., and Webber, B. (2000). Measuring sparticle masses in nonuniversal string inspired models at the LHC. *JHEP*, 0009:004.
- Bolognesi, S., Gao, Y., Gritsan, A. V., Melnikov, K., Schulze, M., et al. (2012). On the spin and parity of a single-produced resonance at the LHC. *Phys.Rev.*, D86:095031.
- Chen, Y., Di Marco, E., Lykken, J., Spiropulu, M., Vega-Morales, R., et al. (2015). 8D likelihood effective Higgs couplings extraction framework in $h \rightarrow 4\ell$. *JHEP*, 1501:125.
- Cowan, G., Cranmer, K., Gross, E., and Vitells, O. (2010). Asymptotic formulae for likelihood-based tests of new physics. *Eur.Phys.J.*, C71:1554.
- Cranmer, K., Lewis, G., Moneta, L., Shibata, A., and Verkerke, W. (2012). HistFactory: A tool for creating statistical models for use with RooFit and RooStats. CERN-OPEN-2012-016, <http://inspirehep.net/record/1236448>.
- Gao, Y., Gritsan, A. V., Guo, Z., Melnikov, K., Schulze, M., et al. (2010). Spin determination of single-produced resonances at hadron colliders. *Phys.Rev.*, D81:075022.
- Gutmann, M. U., Dutta, R., Kaski, S., and Corander, J. (2014). Likelihood-free inference via classification. <http://arxiv.org/abs/1407.4981>.

- Hörmander, L. (1990). *The Analysis of Linear Partial Differential Operators I*. Springer Science, Business Media.
- Ihler, A., Fisher, J., and Willsky, A. (2004). Nonparametric Hypothesis Tests for Statistical Dependency. *IEEE Transactions on Signal Processing*, 52(8):2234–2249. <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=1315943>.
- Jaakkola, T. and Haussler, D. (1998). Exploiting Generative Models in Discriminative Classifiers. <http://citeseer.ist.psu.edu/viewdoc/summary?doi=10.1.1.44.7709>.
- Neal, R. M. (2007). Computing likelihood functions for high-energy physics experiments when distributions are defined by simulators with nuisance parameters. In *Proceedings of PhyStat2007, CERN-2008-001*, pages 111–118. <http://inspirehep.net/record/776337>.
- Ng, A. Y. and Jordan, M. I. (2002). On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes. In Dietterich, T., Becker, S., and Ghahramani, Z., editors, *Advances in Neural Information Processing Systems 14*, pages 841–848. MIT Press.
- Scott, C. and Nowak, R. (2005). A neyman-pearson approach to statistical learning. *IEEE Trans. Inform. Theory*, 51:3806–3819. <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.103.6850&rank=5>.
- Sjostrand, T., Mrenna, S., and Skands, P. Z. (2006). PYTHIA 6.4 Physics and Manual. *JHEP*, 0605:026.
- The ATLAS Collaboration (2012). Observation of a new particle in the search for the

- Standard Model Higgs boson with the ATLAS detector at the LHC. *Phys.Lett.*, B716:1–29.
- The CMS Collaboration (2012). Observation of a new boson at a mass of 125 GeV with the CMS experiment at the LHC. *Phys.Lett.*, B716:30–61.
- Volobouev, I. (2011). Matrix Element Method in HEP: Transfer Functions, Efficiencies, and Likelihood Normalization. <http://arxiv.org/abs/1101.2259>.
- Whiteson, S. and Whiteson, D. (2007). Stochastic optimization for collision selection in high energy physics. *Association for the Advancement of Artificial Intelligence*, 292.
- Xin Tong (2013). A Plug-in Approach to Neyman-Pearson Classification. *Journal of Machine Learning Research*, 14:3011–3040.
- Zadrozny, B. and Elkan, C. (2001). Obtaining calibrated probability estimates from decision trees and naive bayesian classifiers. In *In Proceedings of the Eighteenth International Conference on Machine Learning*, pages 609–616. Morgan Kaufmann. <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.29.3039>.