

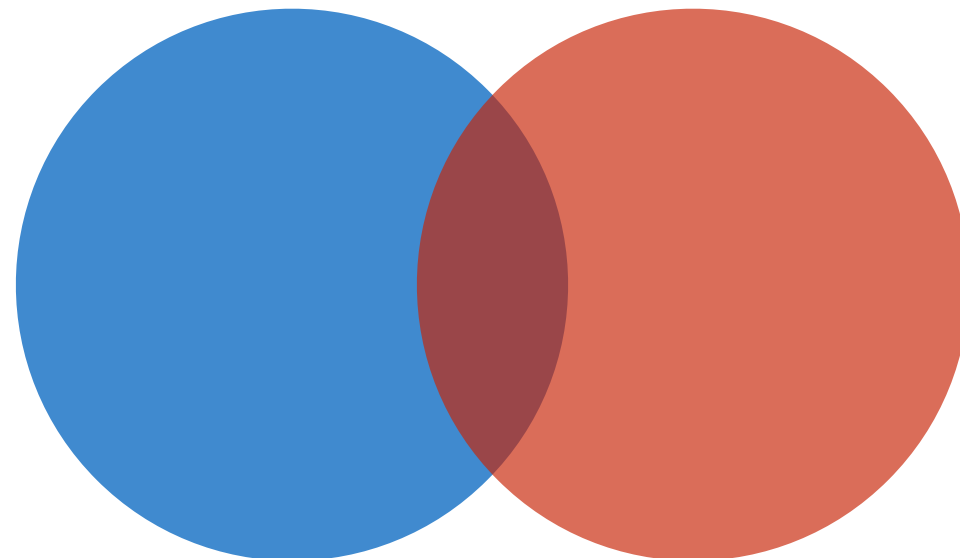
Classification

Goal: Given a *feature vector*, return an integer indexed by the set of possible *classes*.

In most cases, we care about *binary* classification in which there are only two classes (signal versus background)

There are some cases where we care about *multi-class classification*

Feature vector
can be many-
dimensional



Harder = more
overlap between
for **S** and **B**

Classification

Goal: Given a *feature vector*, return an integer indexed by the set of possible *classes*.

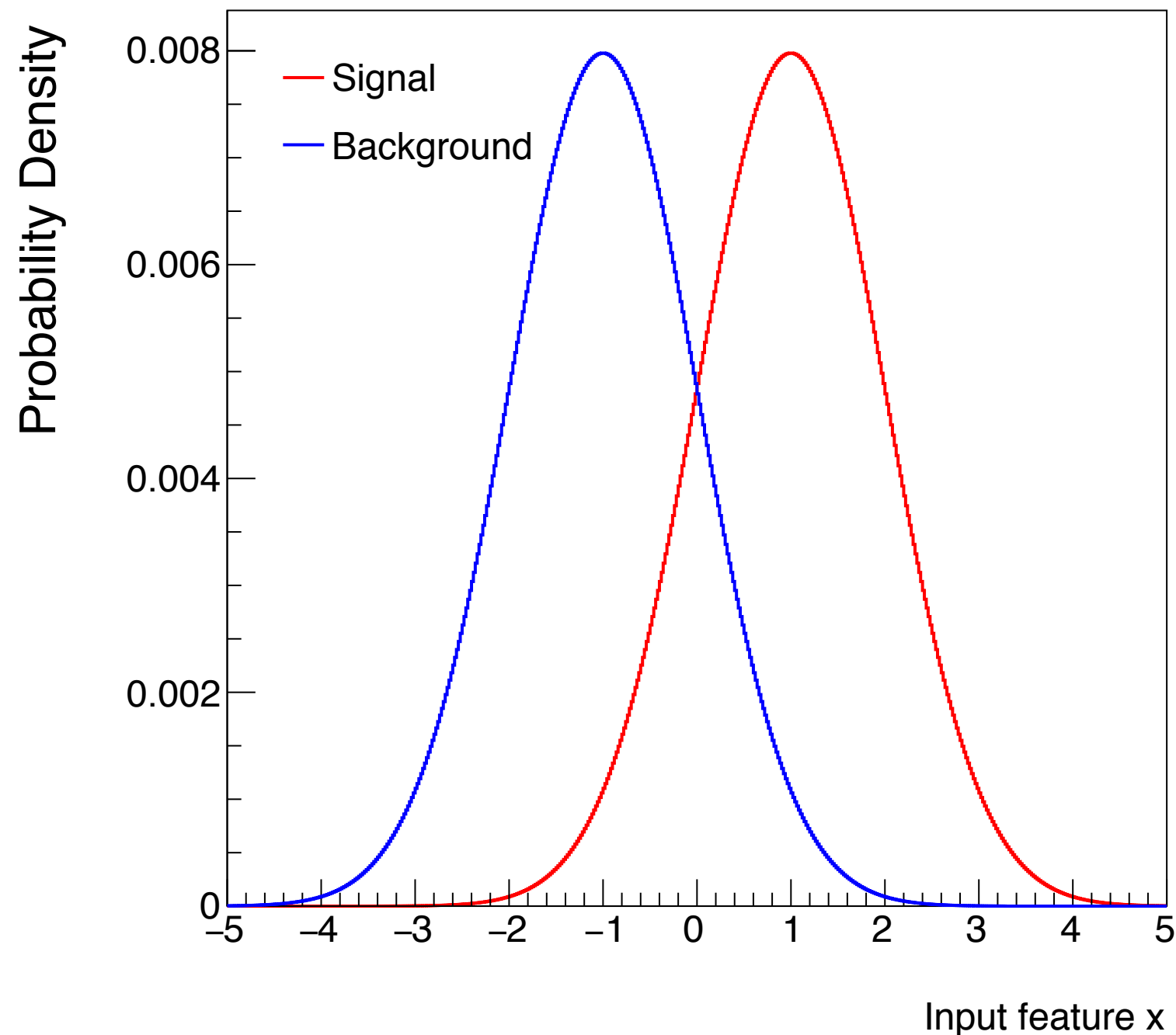
In practice, we don't just want one classifier, but an entire set of classifiers indexed by:

True Positive Rate = signal efficiency =
 $\Pr(\text{label signal} \mid \text{signal}) = \text{sensitivity}$

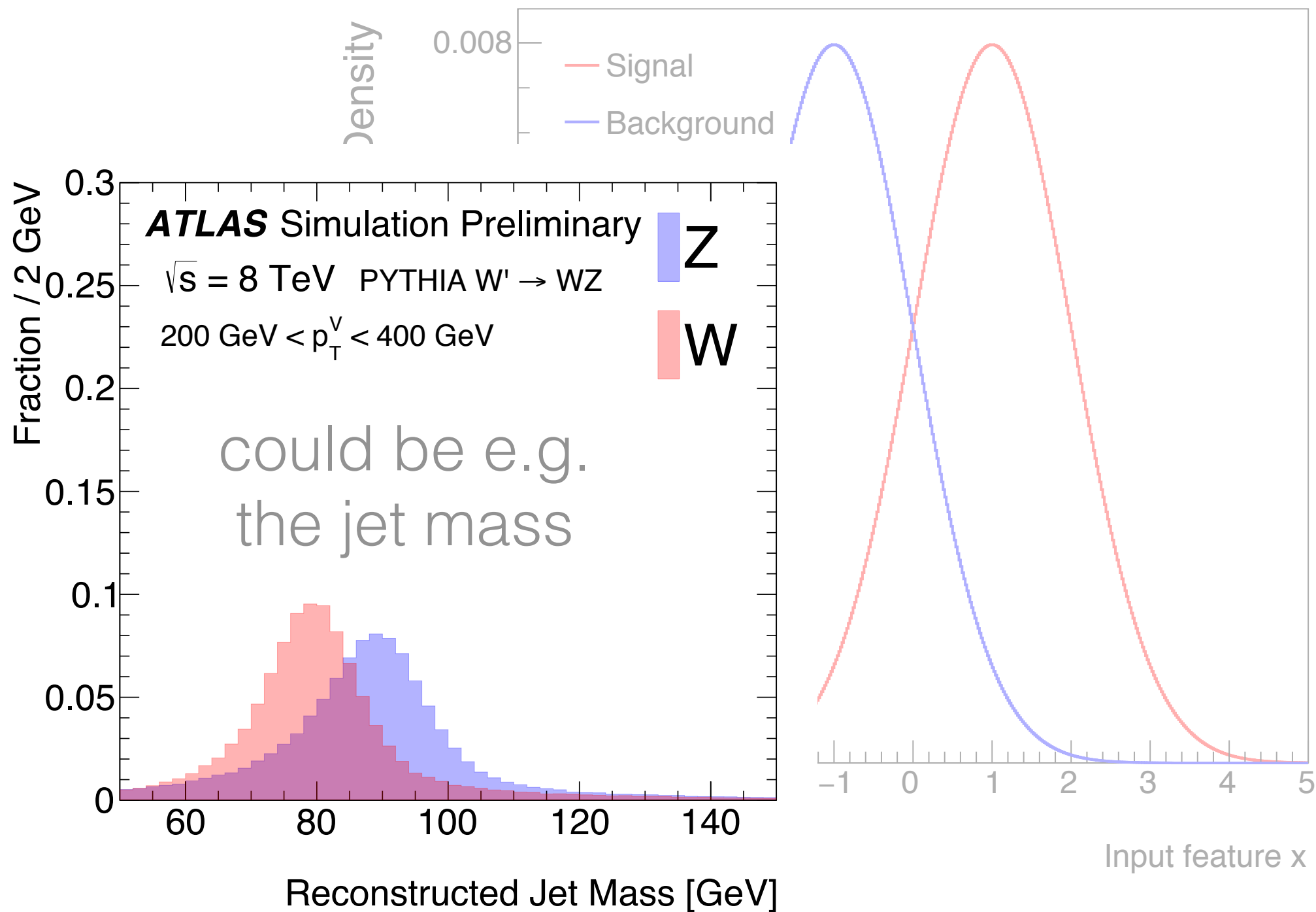
True Negative Rate = 1 - background efficiency =
rejection = $\Pr(\text{label background} \mid \text{background}) = \text{specificity}$

For a given TPR, we want the lowest possible TNR!

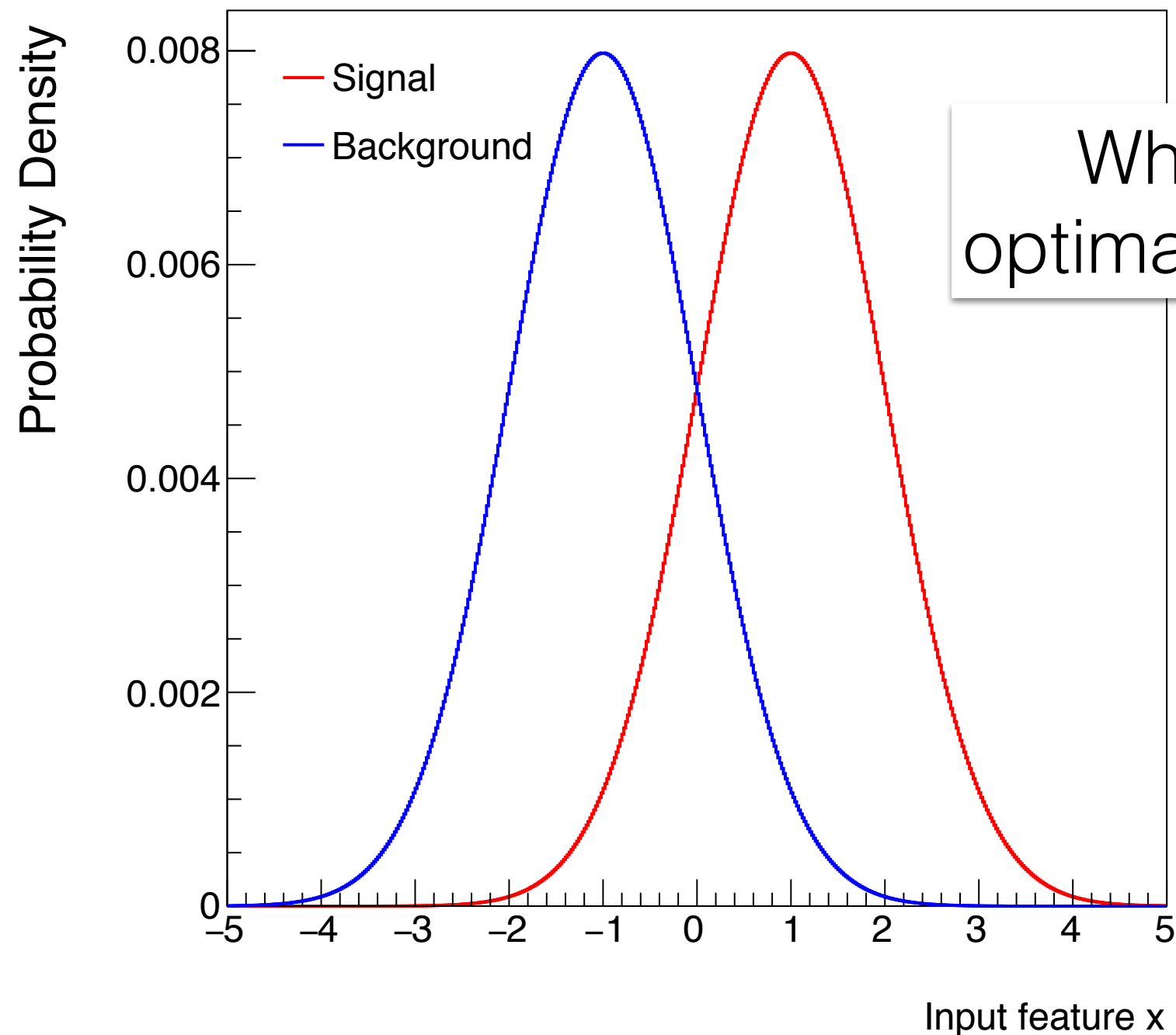
Let's consider an important special case:
binary classification in 1D



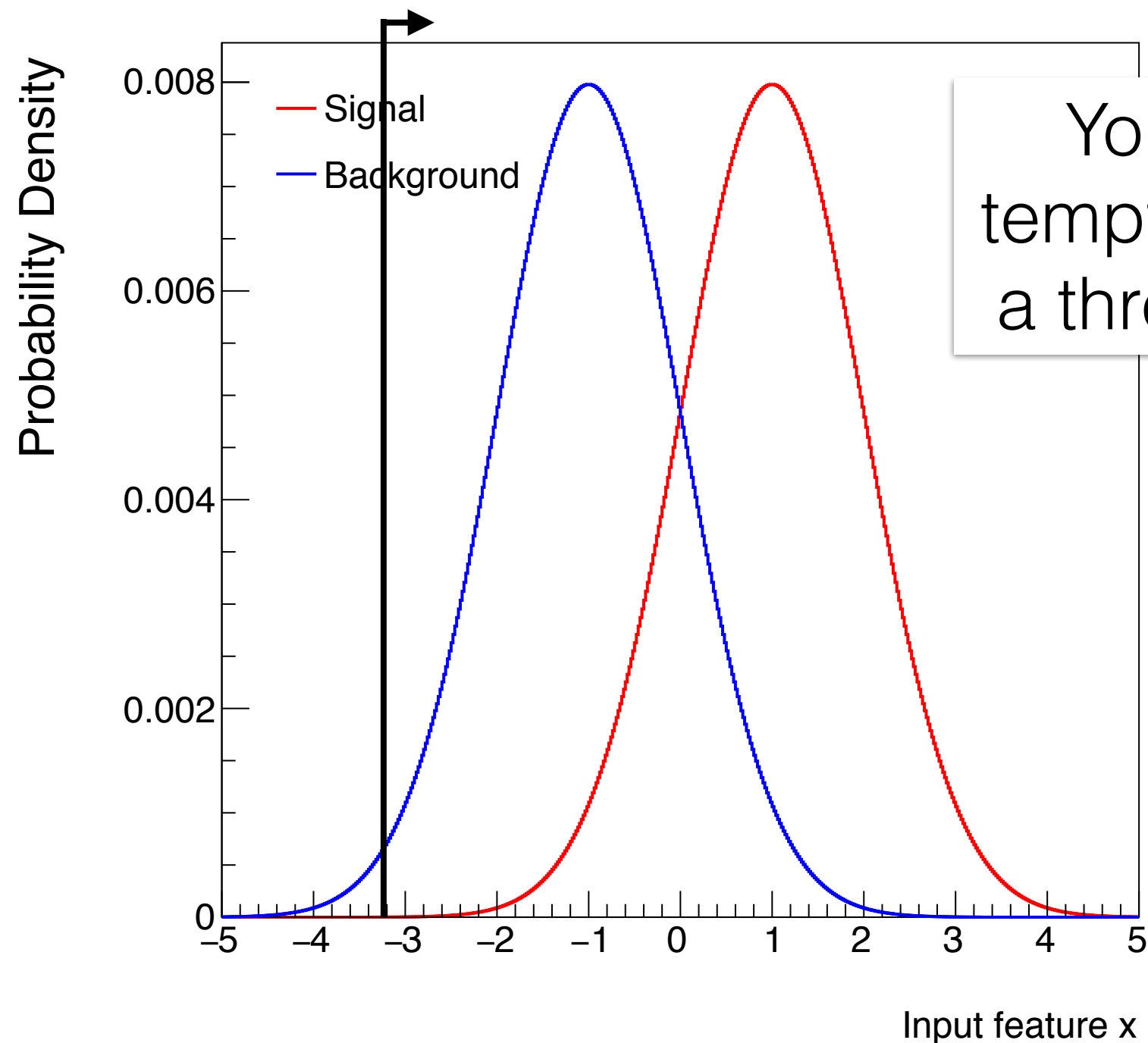
Let's consider an important special case:
binary classification in 1D



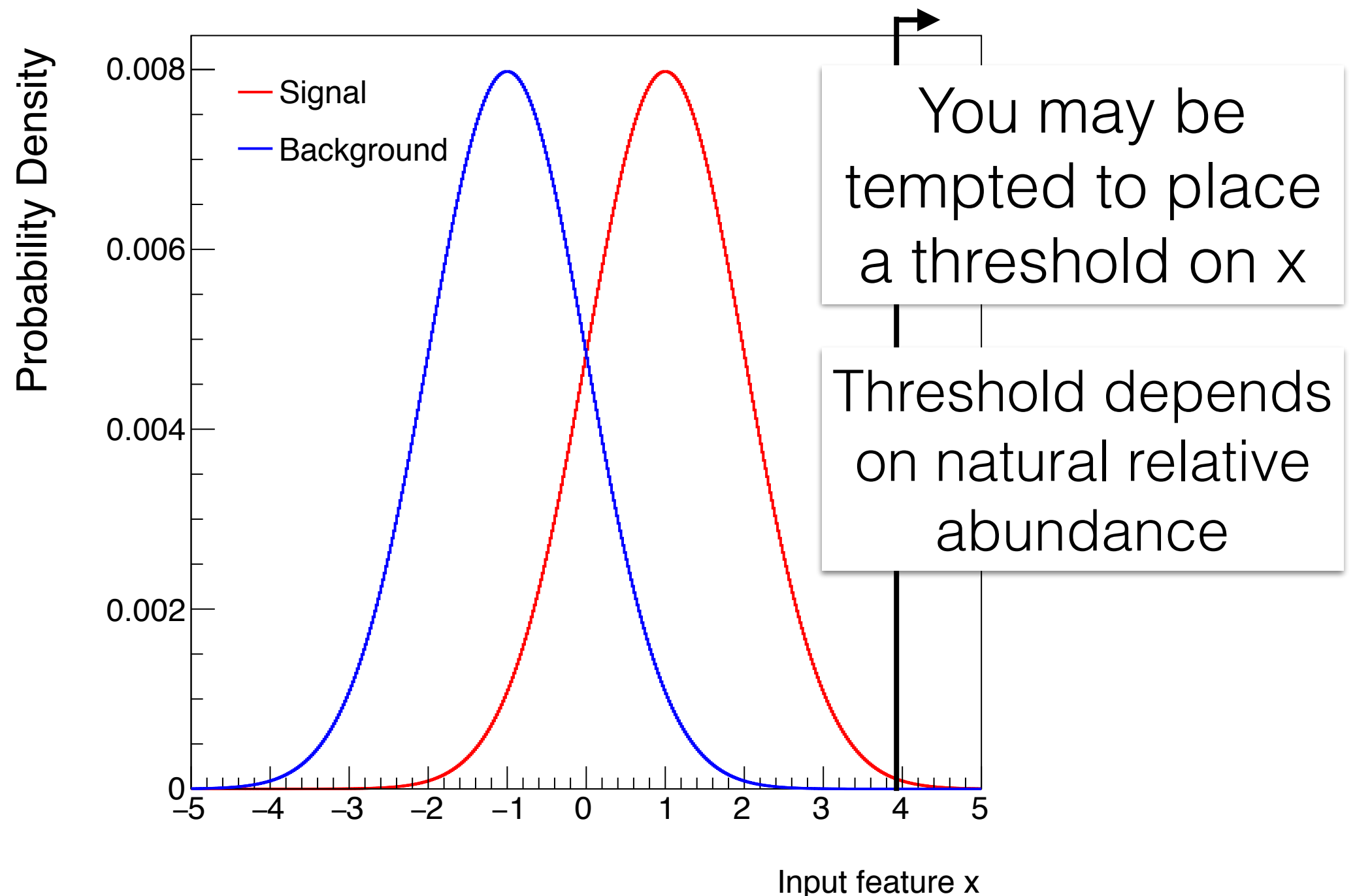
Let's consider an important special case:
binary classification in 1D

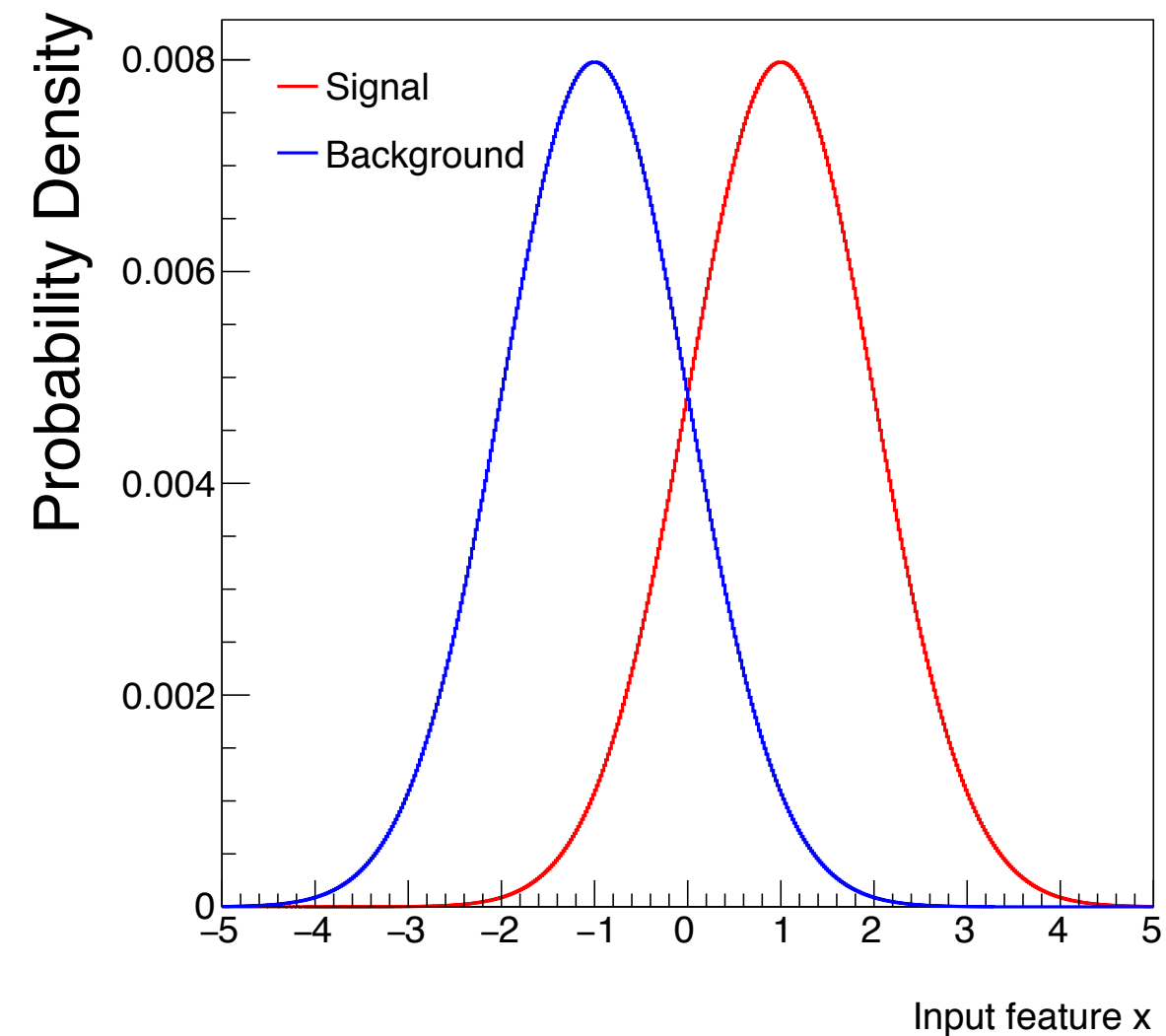


Let's consider an important special case:
binary classification in 1D

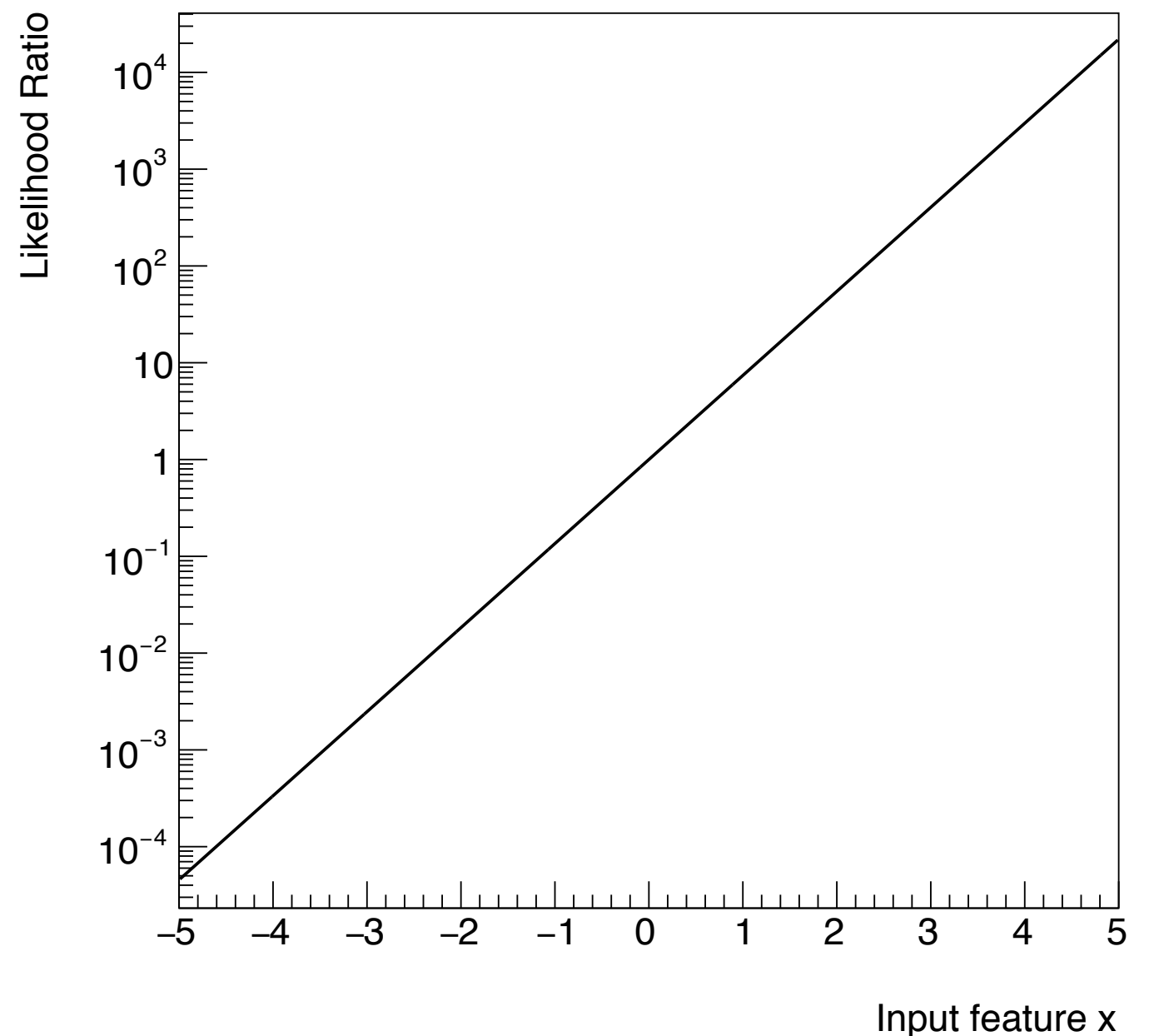


Let's consider an important special case:
binary classification in 1D



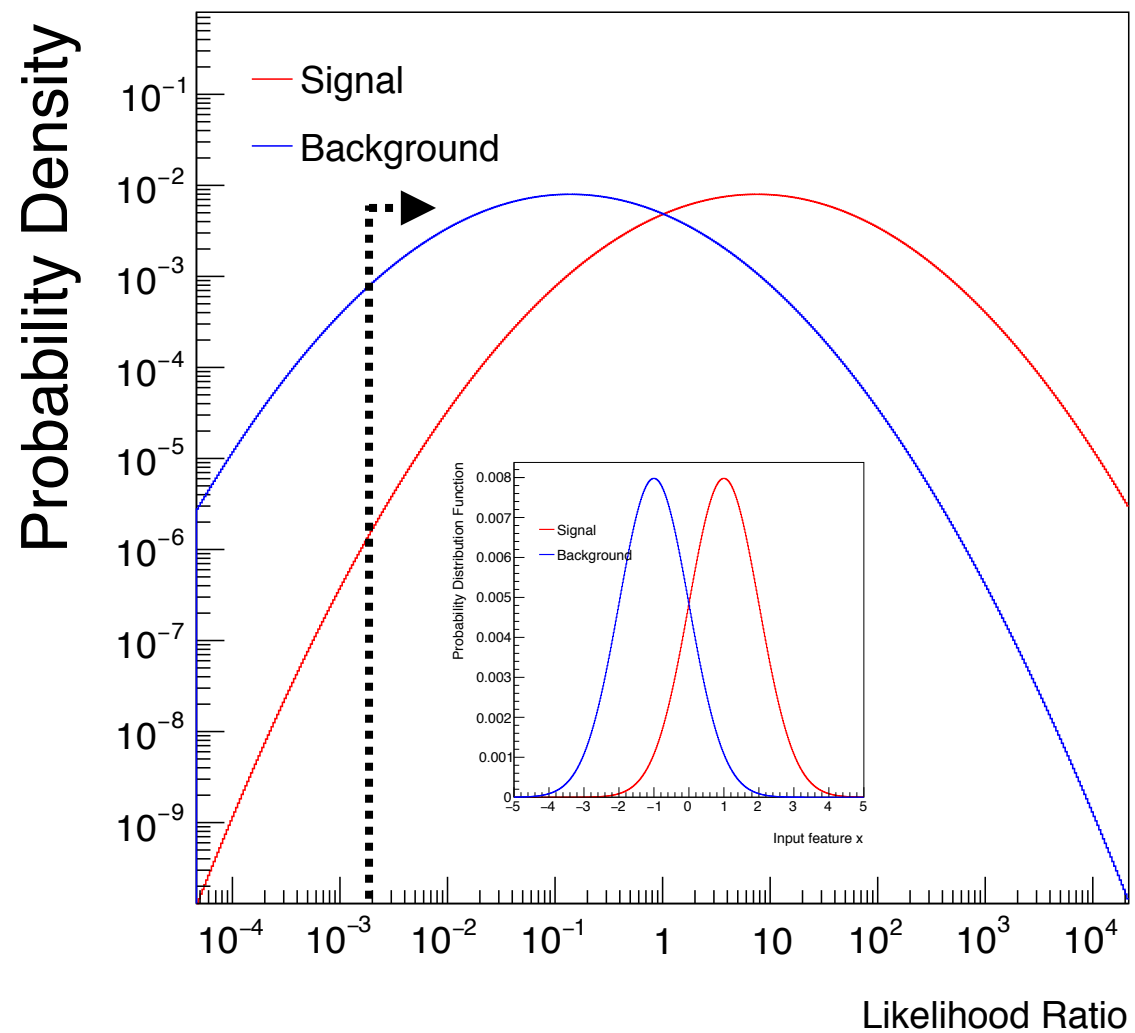


Is the simple threshold cut optimal?



In this simple case, the log LL is proportional to x:
no need for non-linearities!

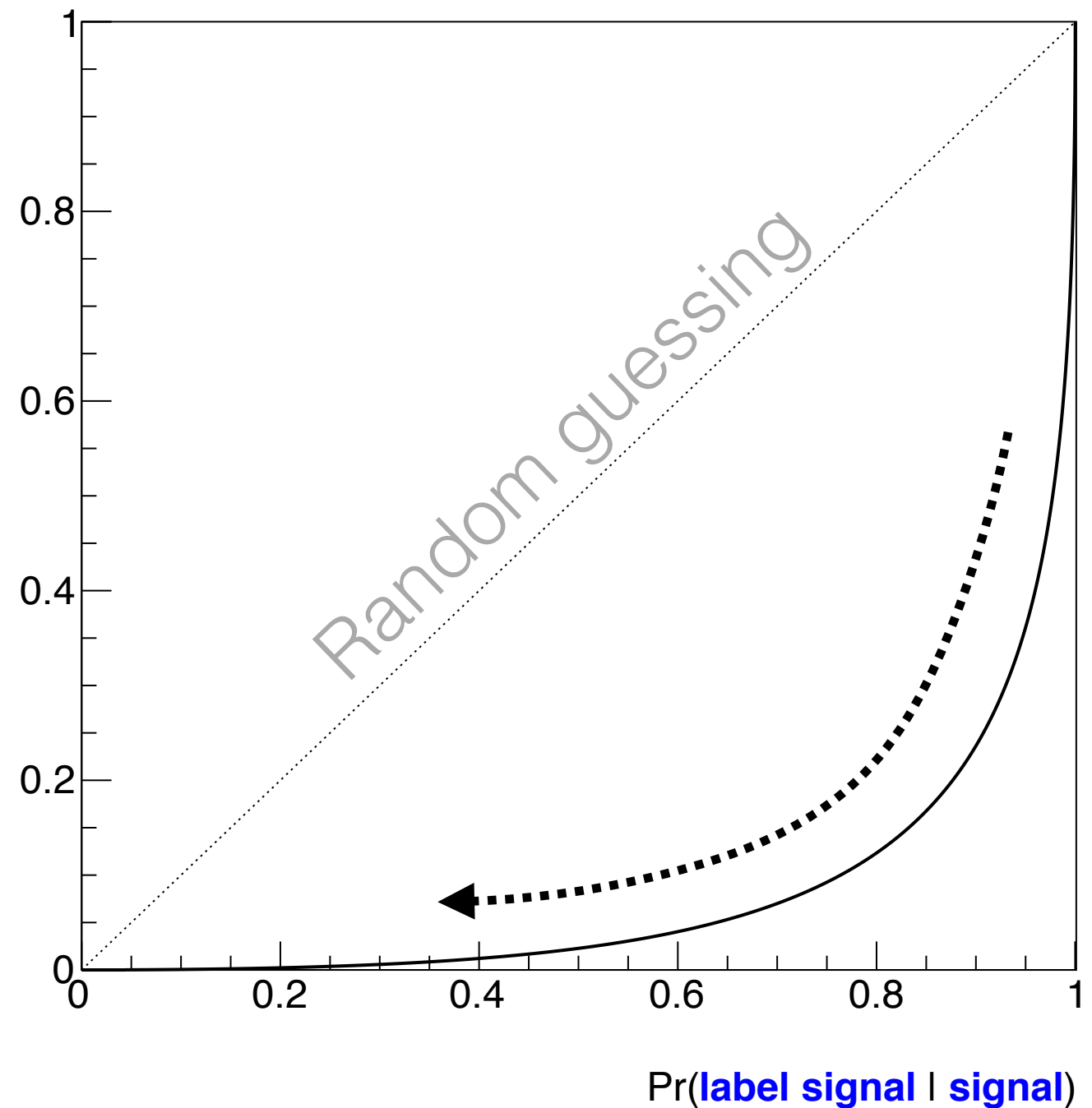
Threshold cut is optimal



“Receiver Operating Characteristic” (**ROC**) Curve

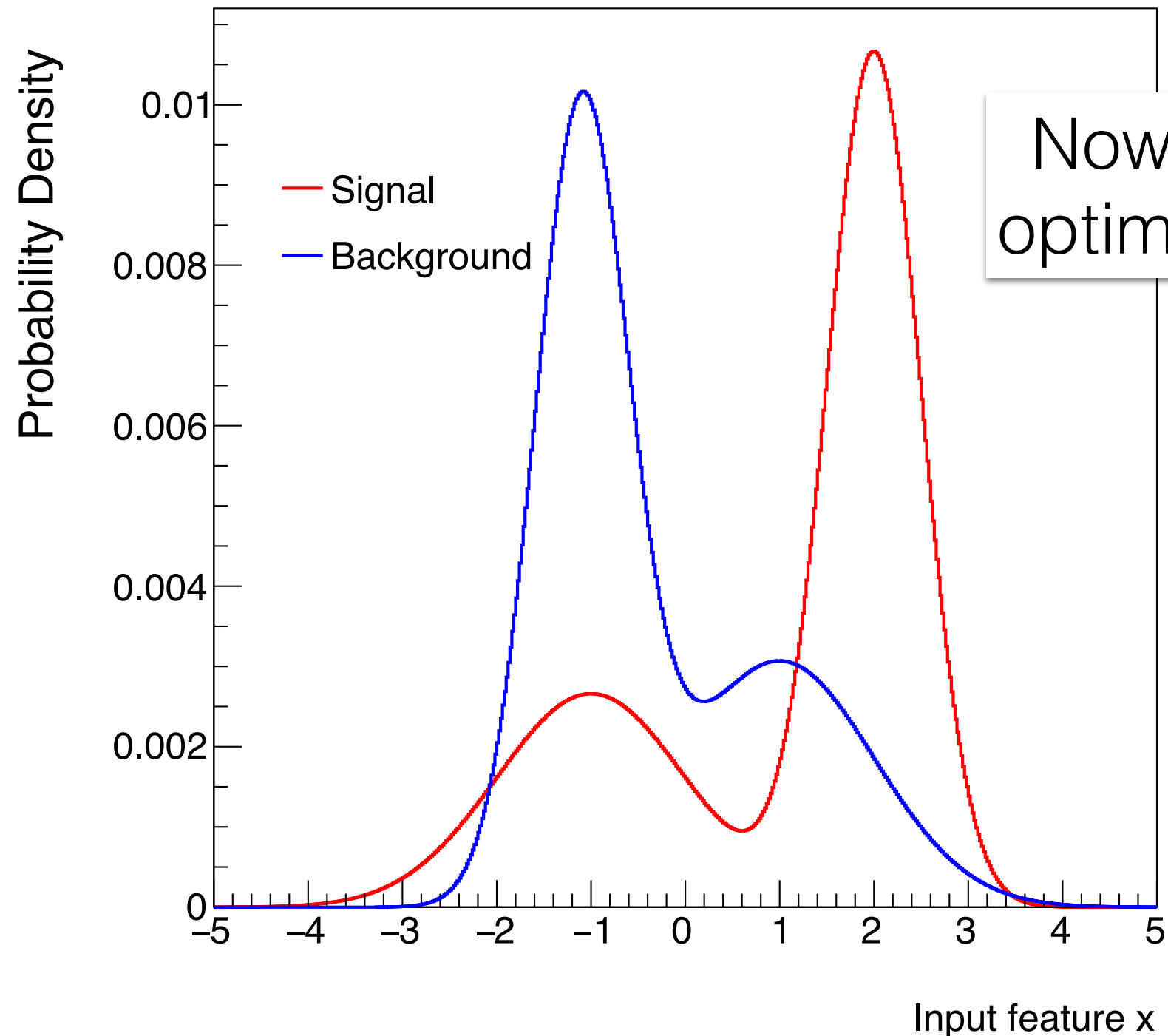
The optimal procedure is a threshold on the LL

$\Pr(\text{label signal} \mid \text{background})$



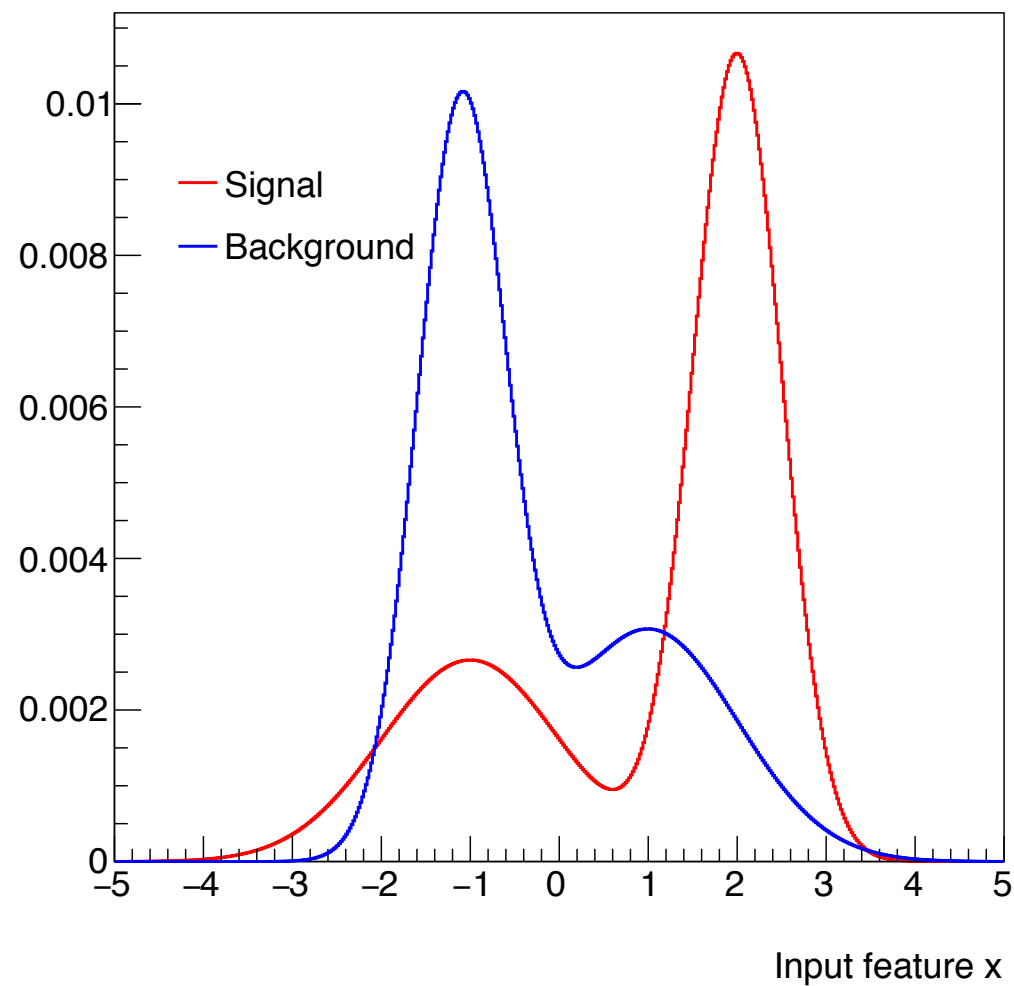
What if the distribution of x is complicated?

Real life is complicated!



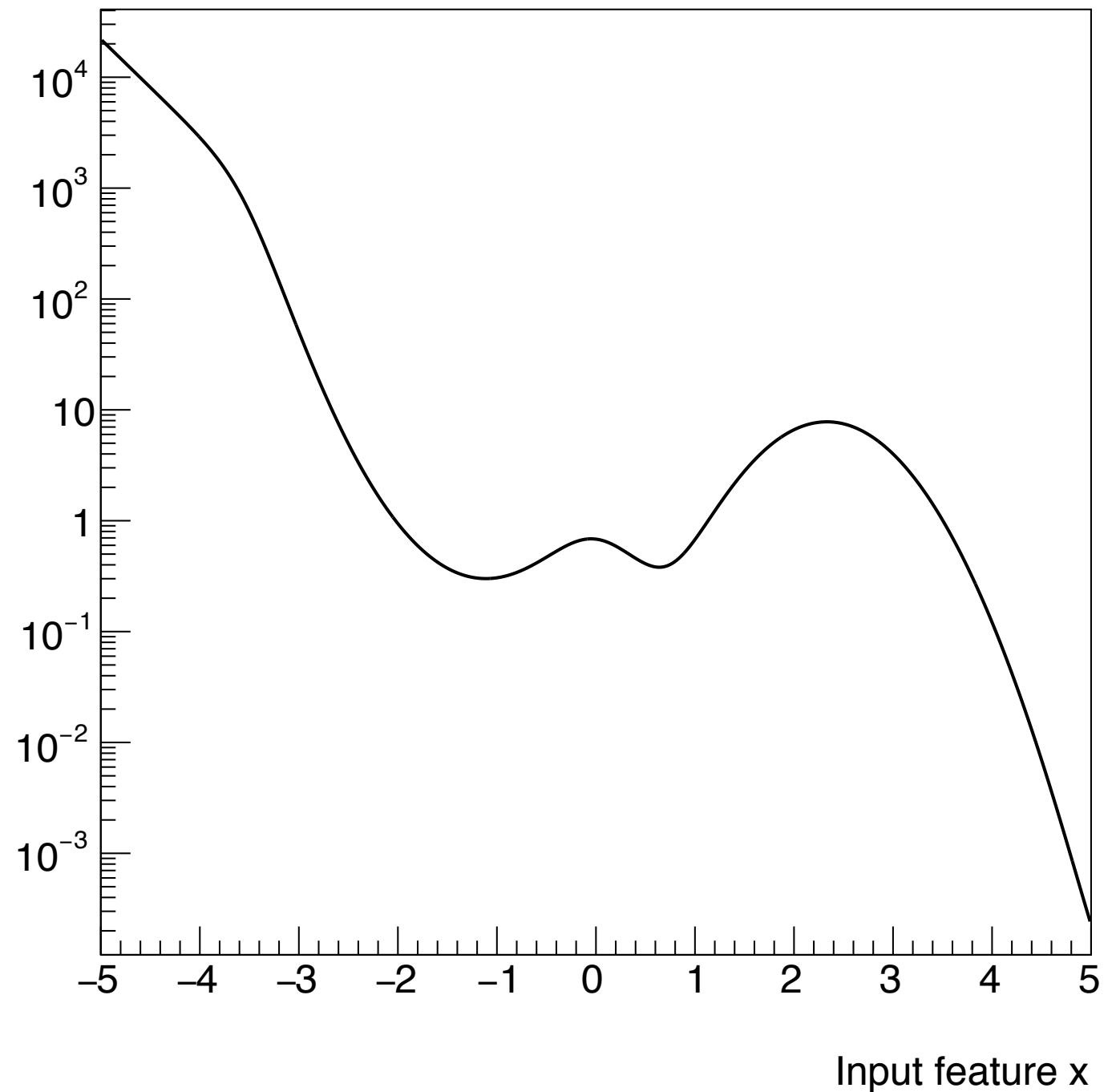
Now what is the optimal classifier?

Probability Density



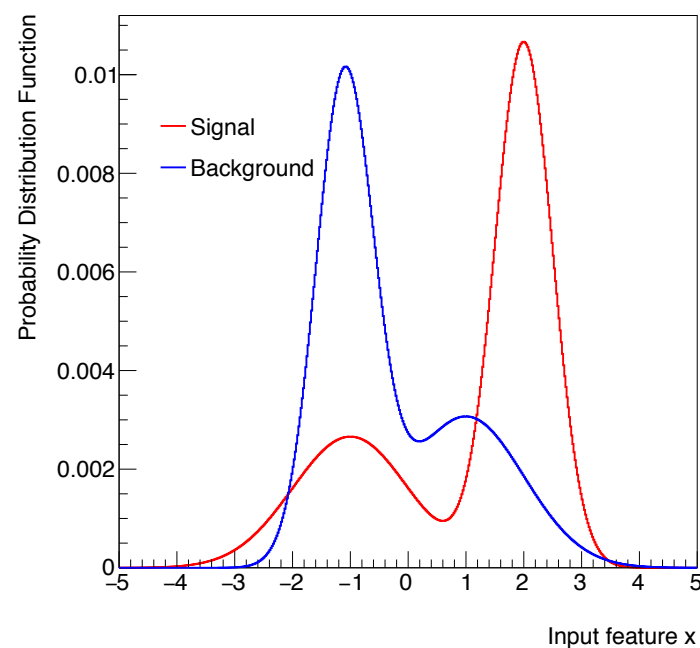
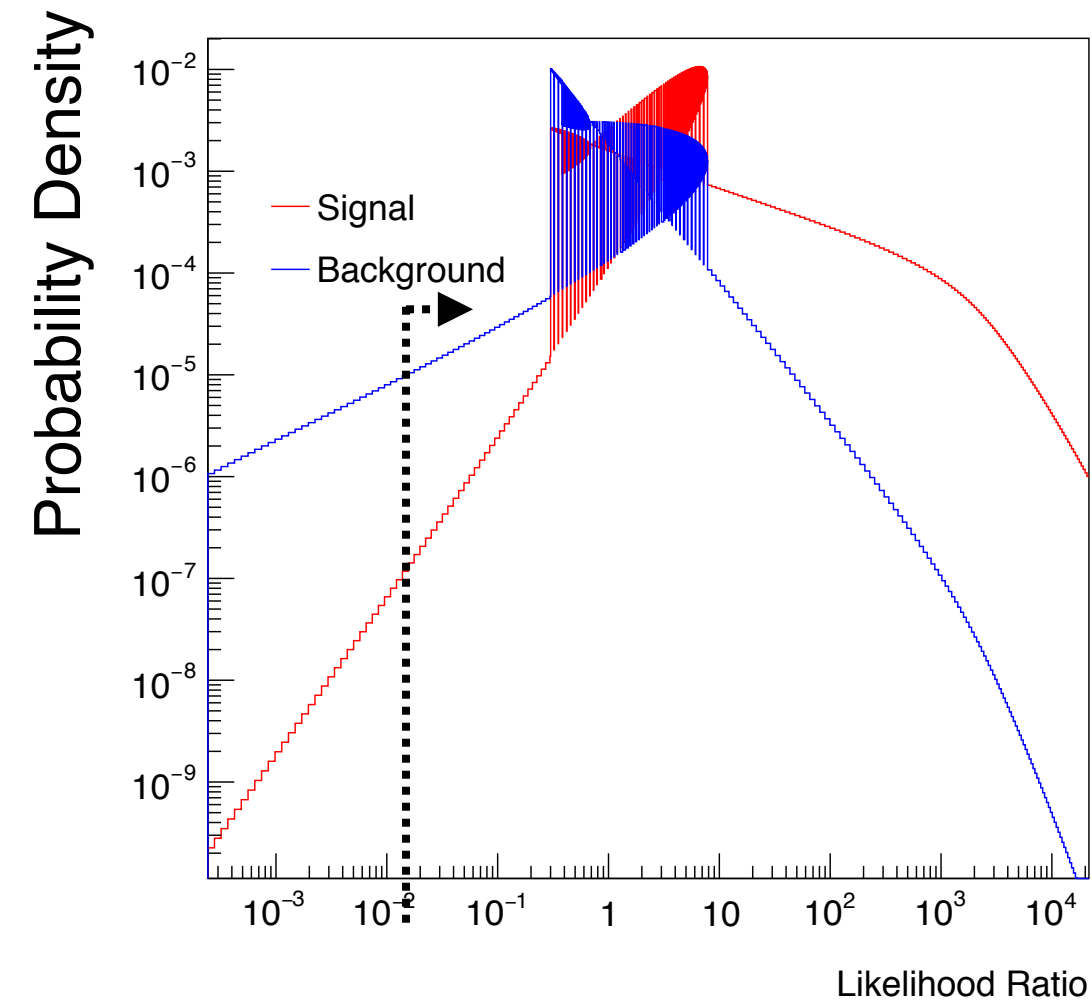
In this case, LL is highly non-linear
(**non-monotonic**) function of x

Likelihood Ratio

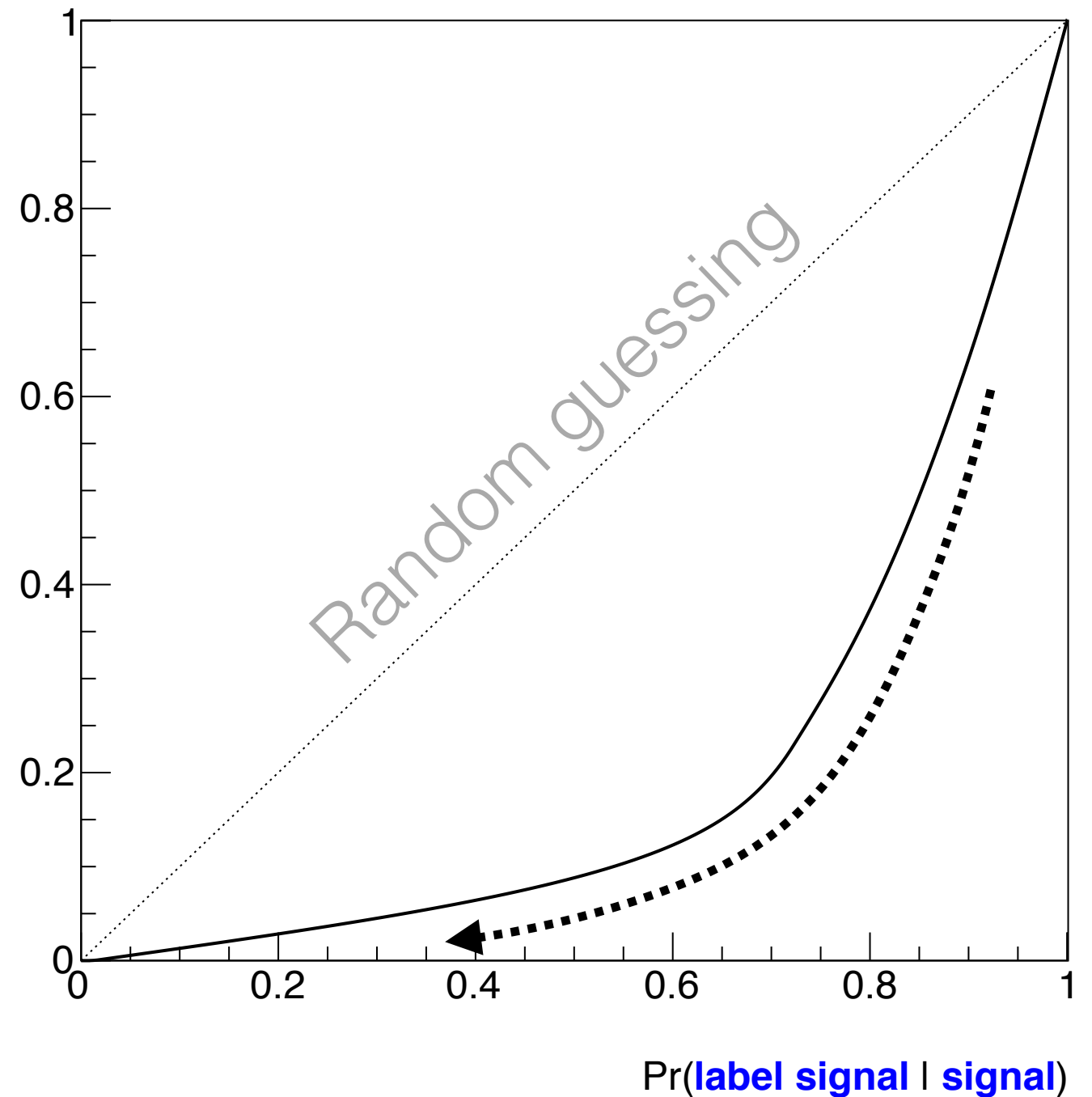


A threshold on x
would be sub-optimal

ROC is worse than the Gaussians,
but that is expected since the
overlap in their PDFs is higher.



$\Pr(\text{label signal} \mid \text{background})$



Why don't we always just
compute the optimal classifier?

In the last slides, we had to estimate the
likelihood ratio - this required binning the PDF

binning works very well in 1D, but becomes
quickly intractable as the feature vector
dimension $\gg 1$ ("curse of dimensionality")

machine learning for classification is simply
**the art of estimating the likelihood ratio
with limited training examples**