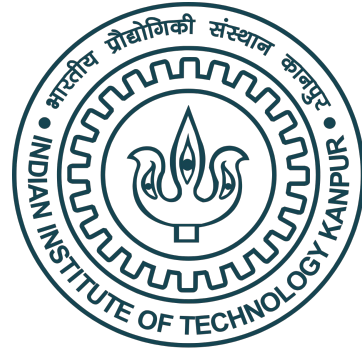


PHY654

Machine learning (ML) in particle physics




Swagata Mukherjee • IIT Kanpur
17th August 2024

Example code for binary classification

NAND Gate using logistic regression

https://github.com/swagata87/IITKanpurPhy654/blob/main/NAND_gate_Logistic_regression.ipynb

NAND



Two features

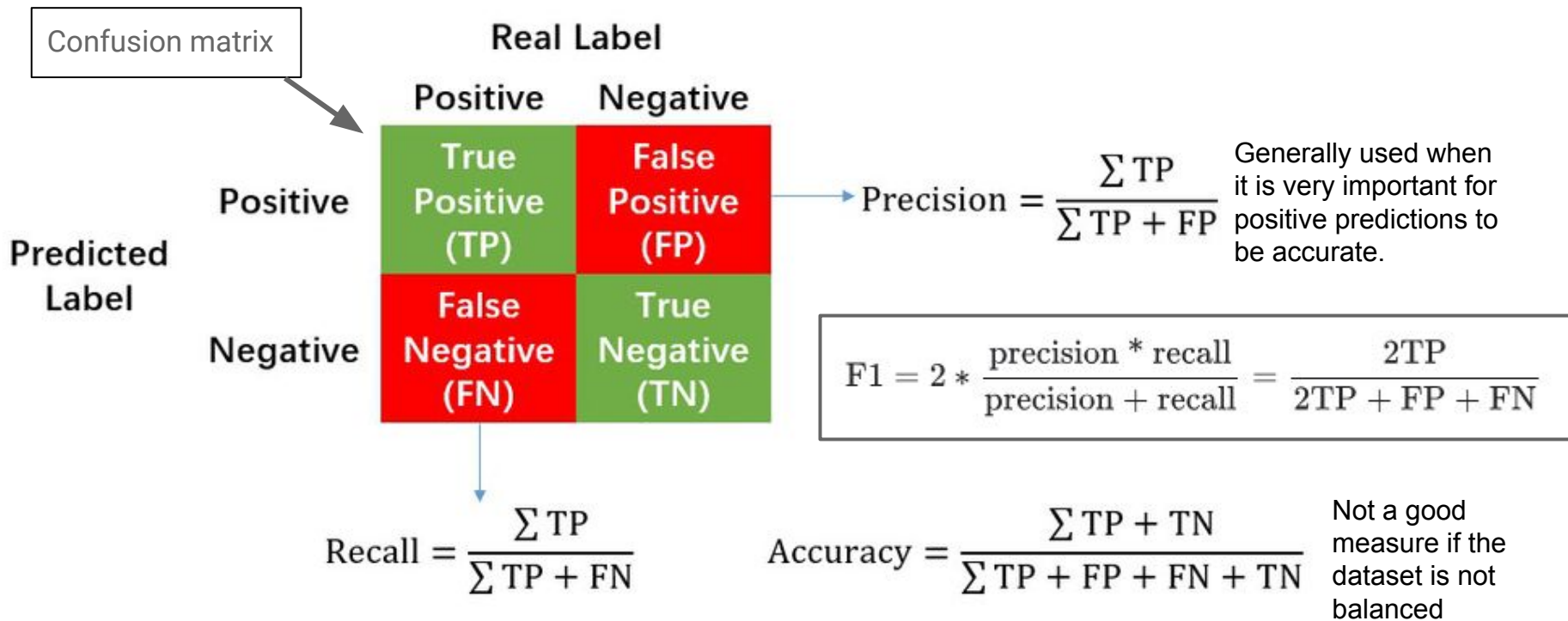
| A | B | Output |
|---|---|--------|
| 0 | 0 | 1 |
| 1 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 1 | 0 |

y



*Try to rewrite the code
without using sklearn library*

**Four training examples
 $m=4$**

Useful metrics for evaluating ML models



An example

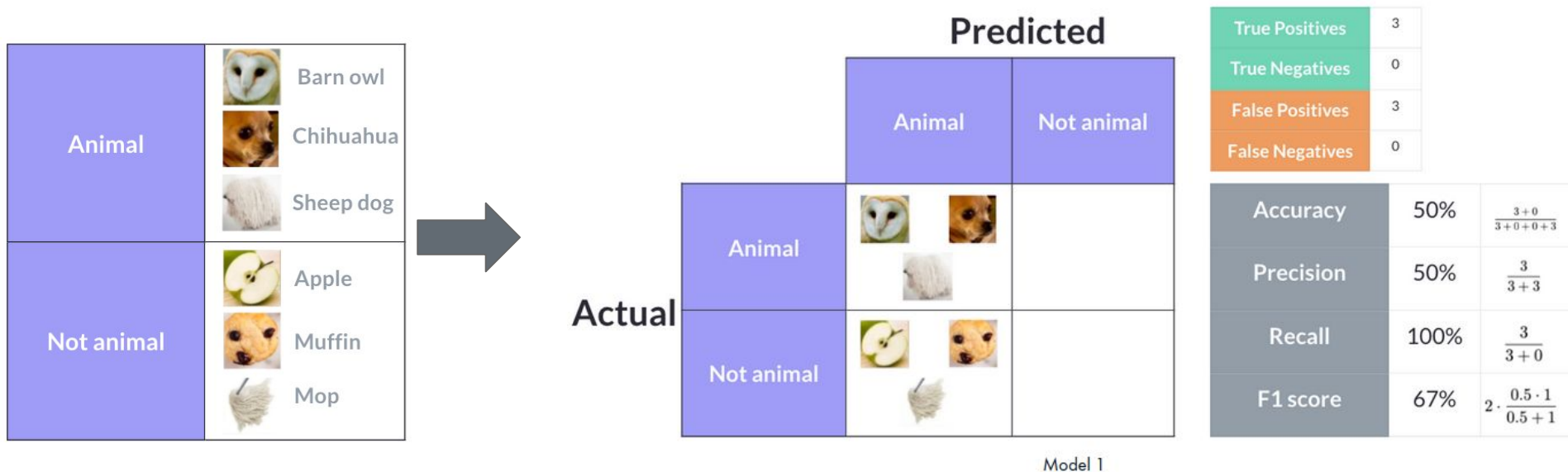
| | | |
|------------|---|-----------|
| Animal |  | Barn owl |
| |  | Chihuahua |
| |  | Sheep dog |
| Not animal |  | Apple |
| |  | Muffin |
| |  | Mop |



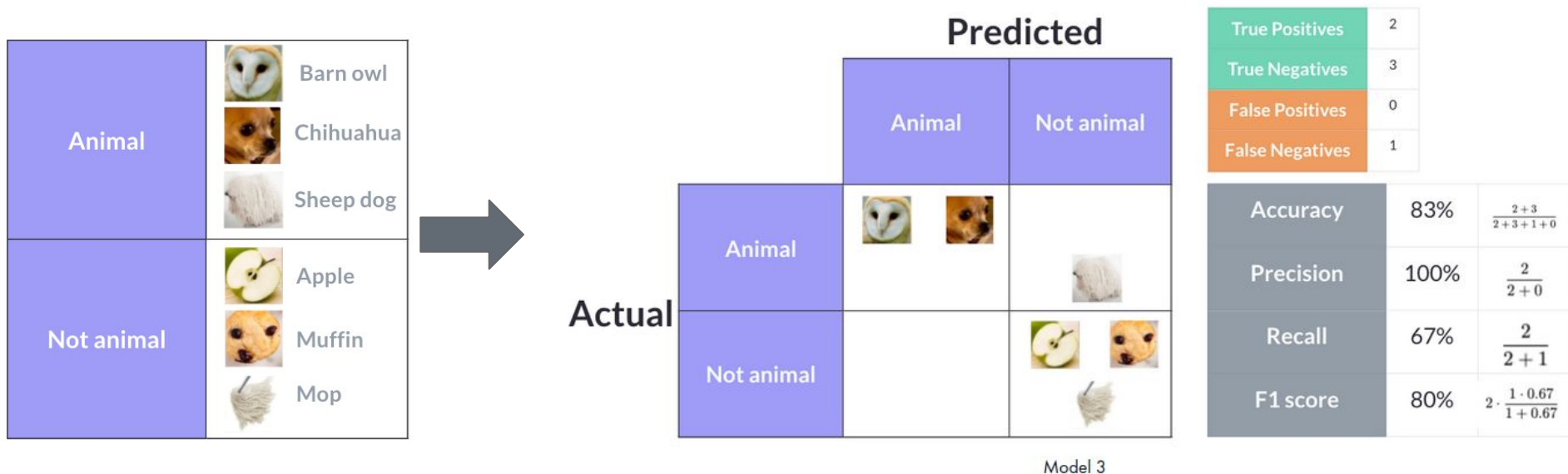
| | | Predicted | |
|--------|------------|-----------------|-----------------|
| | | Animal | Not animal |
| Actual | Animal | True Positives | False Negatives |
| | Not animal | False Positives | True Negatives |

Confusion matrix

Look at all these metrics



Look at all these metrics



These metrics are calculated at a single classification threshold value (predicted output > 0.5 or not). But if you want to evaluate a model's quality across all possible thresholds, you need different tools.

ROC

Receiver-operating characteristic curve (ROC)

ROC curve is drawn by computing true positive rate (TPR) and false positive rate (FPR) at every possible threshold (in selected intervals), then graphing TPR vs FPR.

False Positive Rate = False Positives / (False Positives + True Negatives)

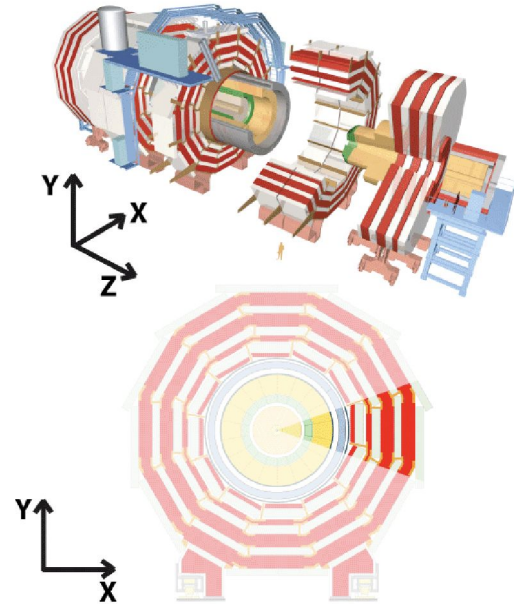
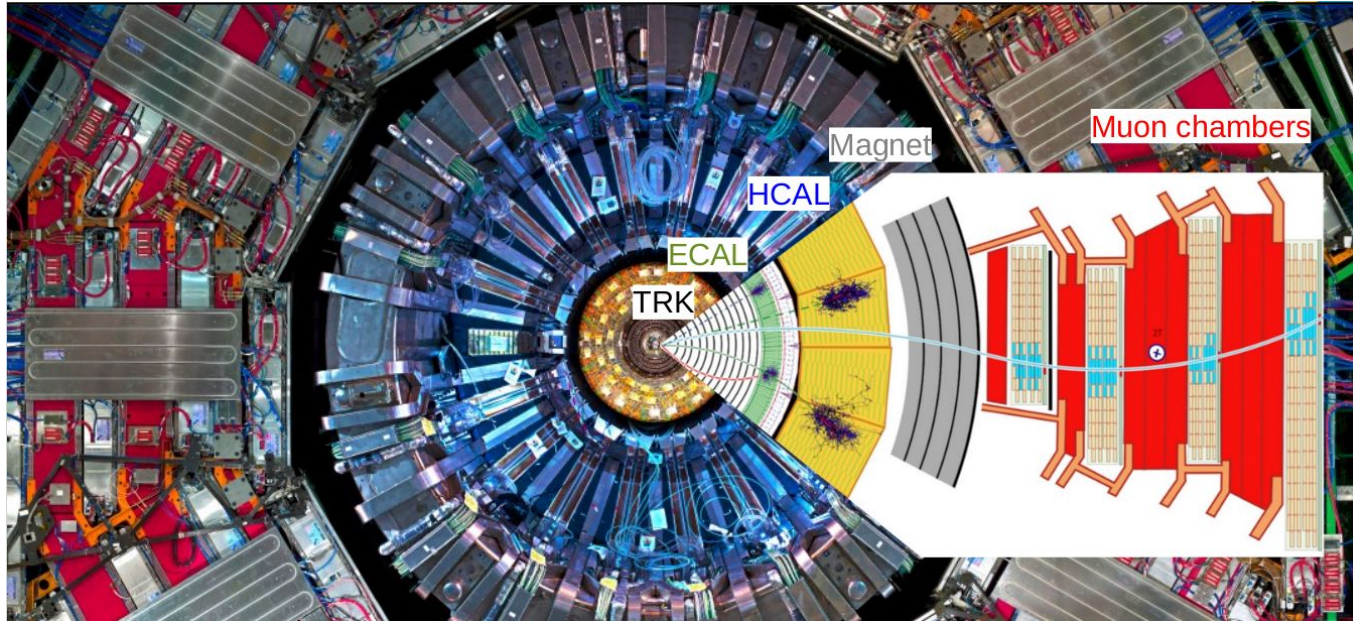
True Positive Rate = True Positives / (True Positives + False Negatives)

In particle physics, we generally plot **signal efficiency vs background efficiency**; or signal efficiency vs background rejection.

ROC is not specific to ML. It can be plotted for cut-based physics analysis as well. It tells us the usefulness of a cut applied on a variable.

The area under the curve (AUC) can be used as a single-valued metric.

Binary classification in physics experiments

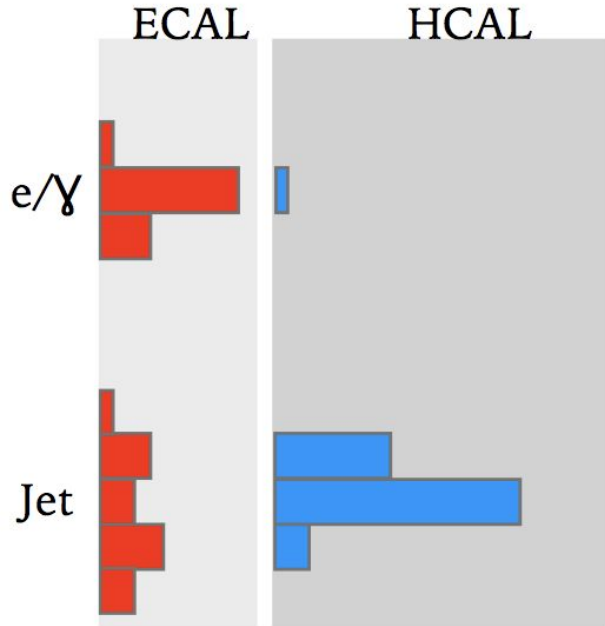


Photon vs hadronic-jet : Binary classification

Photon is signal ($y=1$) and Hadronic-jet is background ($y=0$)

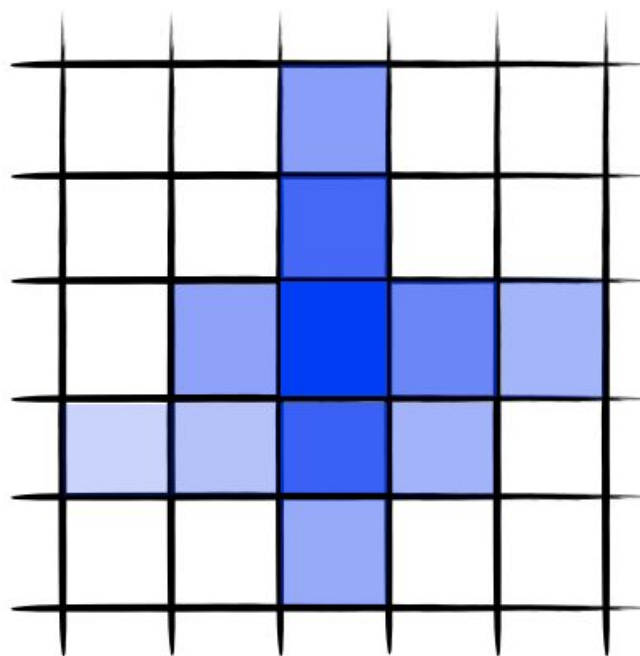
Feature 1 : H/E

Photon vs hadronic-jet : Binary classification



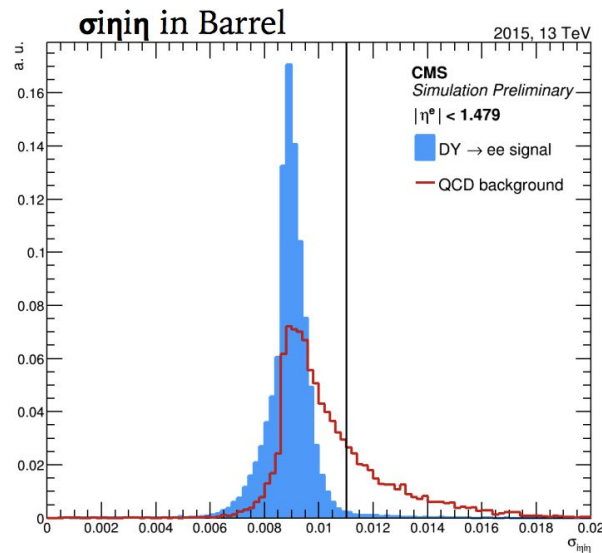
H/E variable is a good discriminator variable.

Feature 2 : Showershape in ECAL



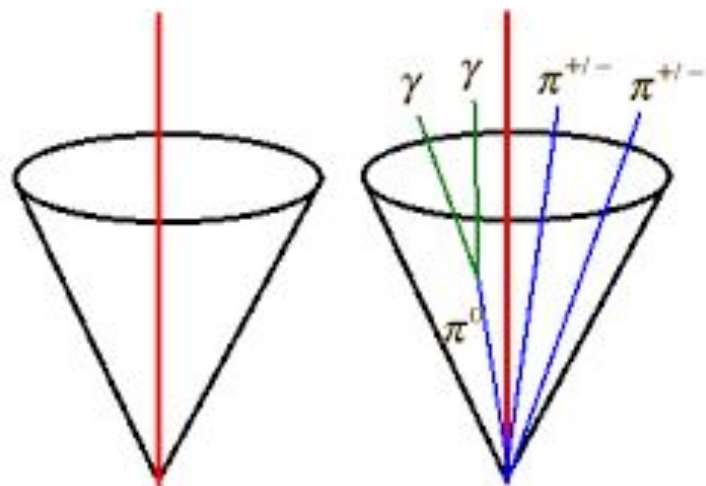
Photon vs hadronic-jet : Binary classification

Showershape is also a good discriminator between signal and background



Feature 3 : Isolation sum

Photon vs hadronic-jet : Binary classification

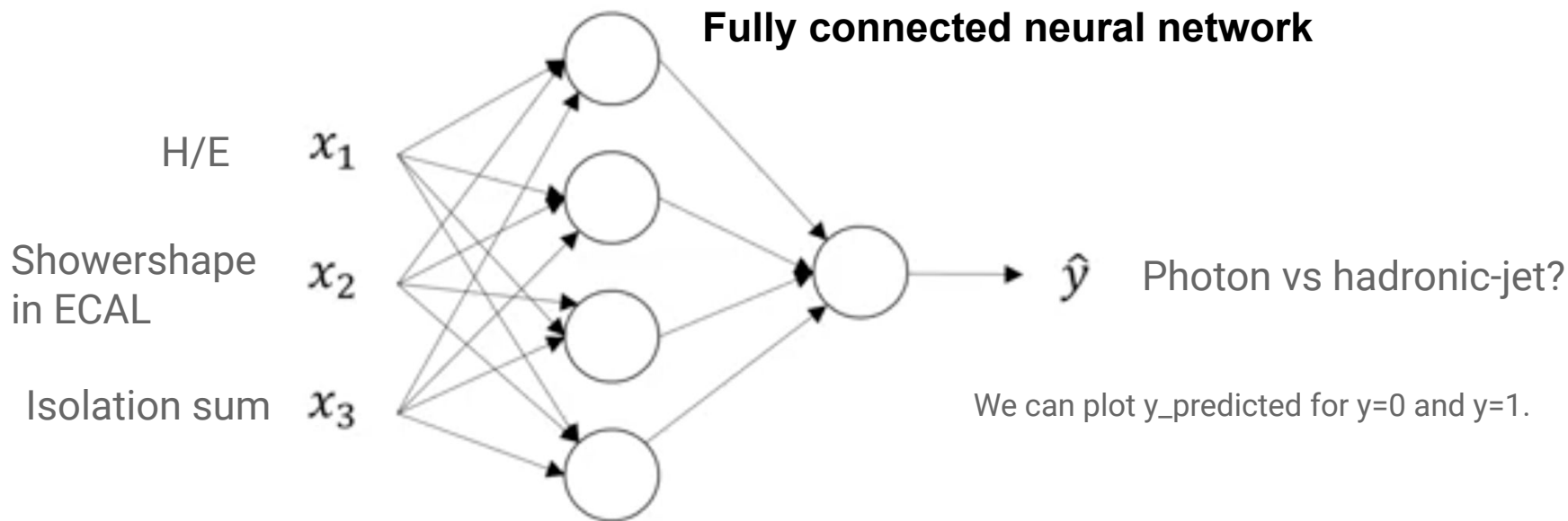


Isolated

Non-Isolated

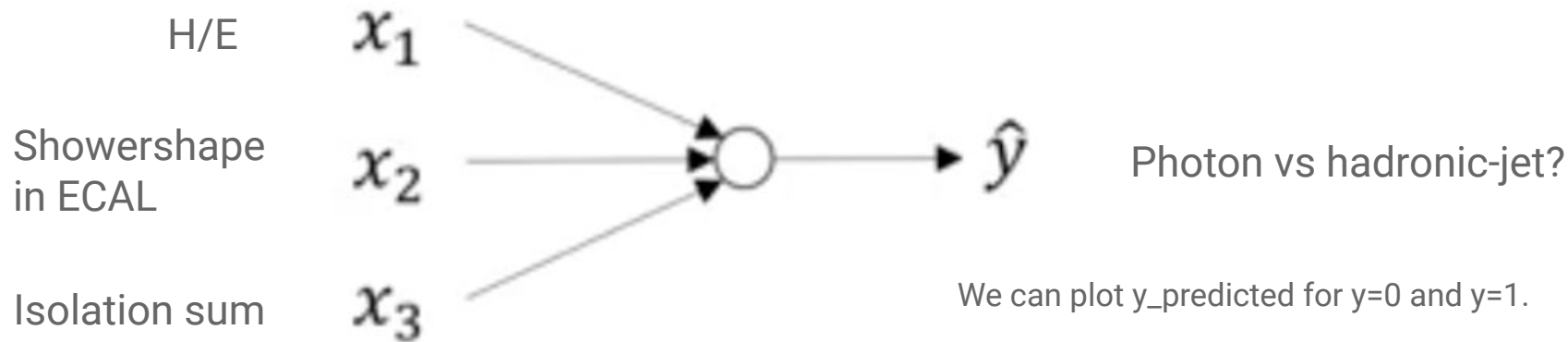
Isolation sum is also a good discriminator between signal and background

Example use-case in physics experiments

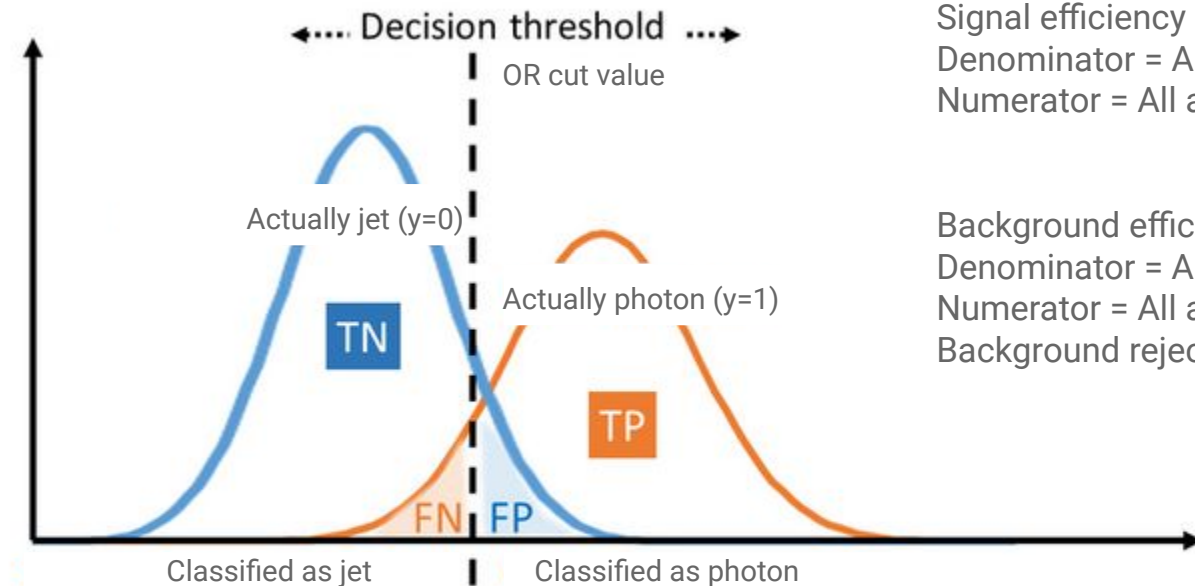


Example use-case in physics experiments

Logistic regression network



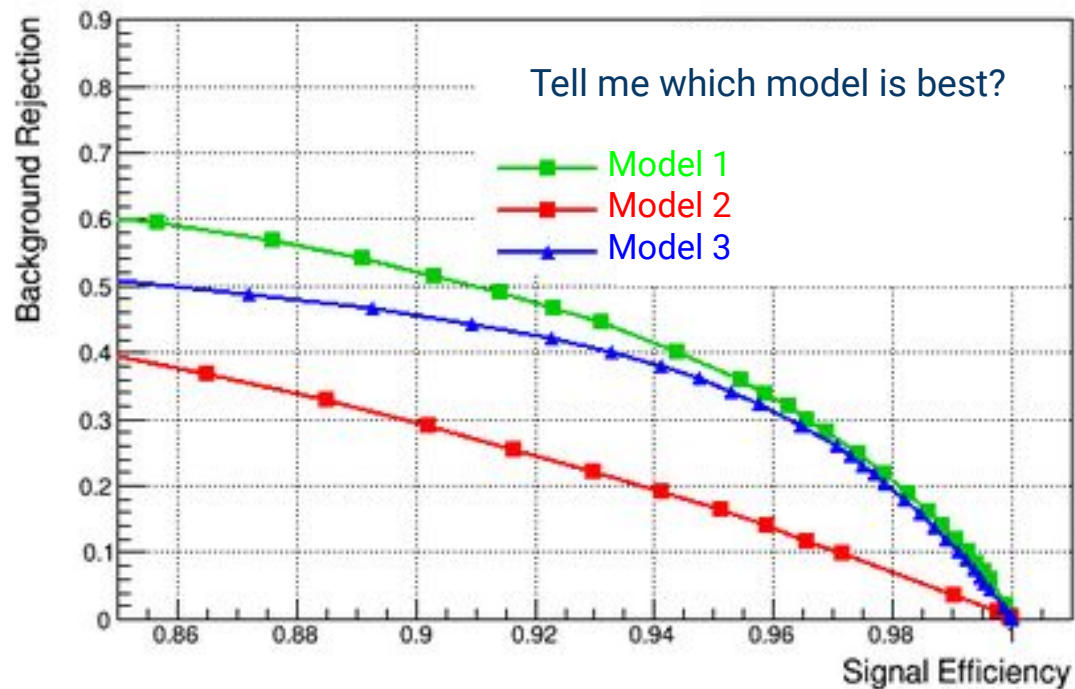
Distribution of predicted output



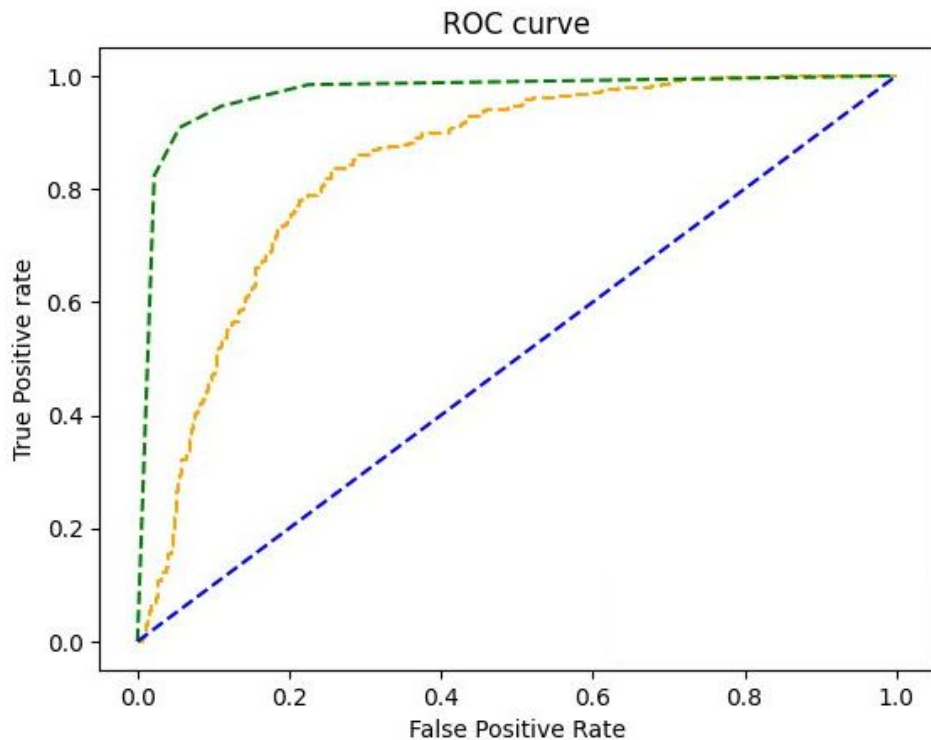
Signal efficiency = numerator / denominator
Denominator = All actual photons
Numerator = All actual photons that pass the cut

Background efficiency = numerator / denominator
Denominator = All actual jets
Numerator = All actual jets that pass the cut
Background rejection = 1 - background efficiency

ROC: an example



ROC: another example



Tell me which model is best?

Some useful links

- Python Tutorial (<https://docs.python.org/3.7/tutorial/index.html>): an introduction to the Python programming language
- Google Colab (<https://colab.research.google.com/>): for Python development in your web-browser
- numpy (<https://numpy.org/doc/stable/user/quickstart.html>): a widely used library for mathematical operations in Python
- Keras (<https://keras.io/>): a beginner-friendly deep learning library
- Tensor Flow (<https://www.tensorflow.org/>): a useful backend for deep learning development
- SciKit Learn (<https://scikit-learn.org/stable/>): helpful machine learning library
- Seaborn (<https://seaborn.pydata.org/>): a library for creating graphs and figures

A non-physics example of classification



→ Input image →

Binary classifier
1 (cat) vs 0 (non cat)

What kind of data we will be dealing with?

Structured data

| H/E | isolation | shower-shape | Matched tracks? | Photon or jet |
|-----|-----------|--------------|-----------------|---------------|
| - | - | - | - | - |
| - | - | - | - | - |
| - | - | - | - | - |

Unstructured data



Image



Audio

Images as input to a Neural Network?



→ Input image →

Binary classifier
1 (cat) vs 0 (non cat)

How does a computer “**see**” an image?

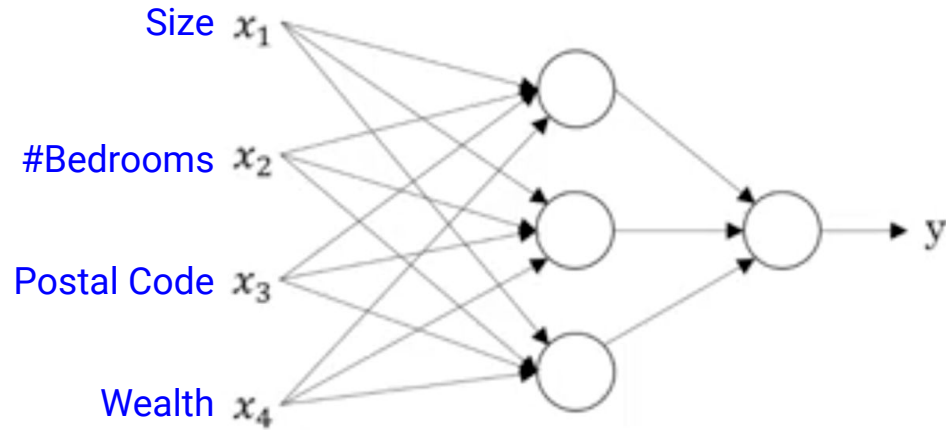
→ It sees 3 matrices. These are pixel-intensity values.

| | | | | | | |
|-------|-----|------|-----|-----|-----|-----|
| | | Blue | | | | |
| Green | | 255 | 134 | 93 | 22 | |
| Red | | 255 | 134 | 202 | 22 | 2 |
| | 255 | 231 | 42 | 22 | 4 | 30 |
| | 123 | 94 | 83 | 2 | 192 | 124 |
| | 34 | 44 | 187 | 92 | 34 | 142 |
| | 34 | 76 | 232 | 124 | 94 | |
| | 67 | 83 | 194 | 202 | | |

→ Each element is a feature. We need to unroll to a column vector.

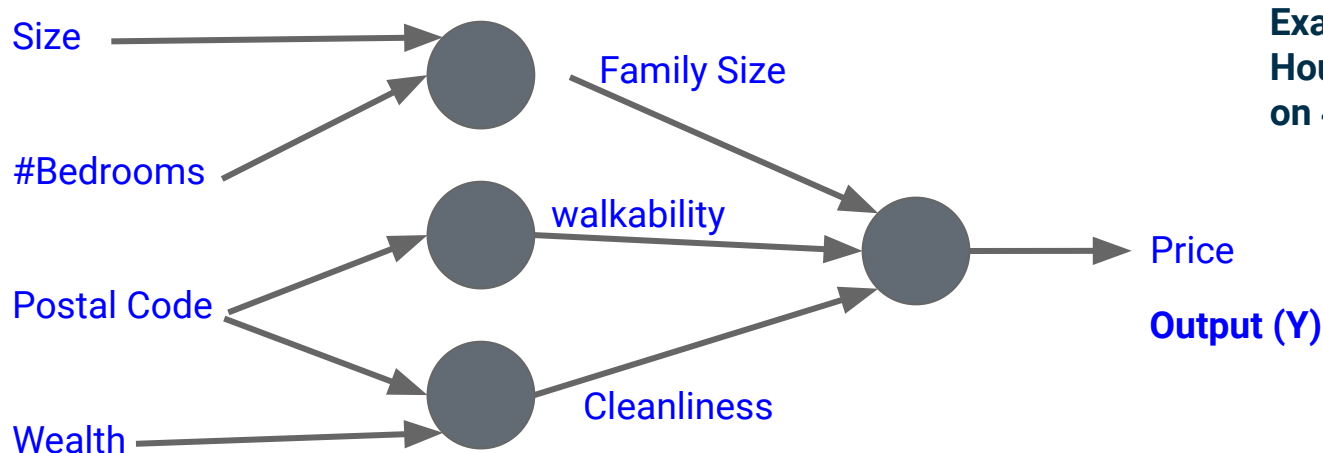
If input image is 64x64 then number of features = $64 \times 64 \times 3 = 12288$.
So, $n=12288$.

Hidden layer: what does it do?



Example (regression):
Housing price prediction, based
on 4 input variables.

Hidden layer: what does it do?



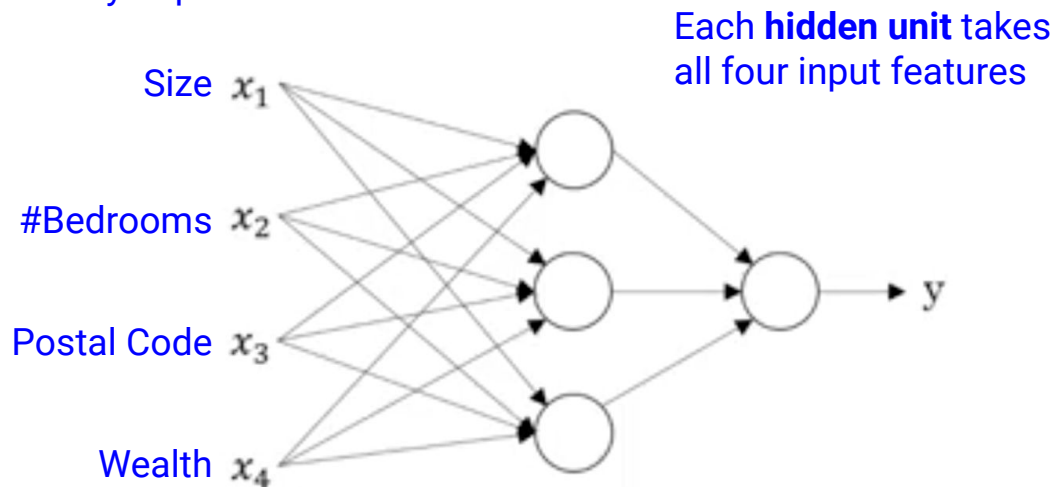
Example (regression):
Housing price prediction, based
on 4 input variables.

Input features (X)

You will provide X and Y for some training examples, and the hidden layer(s) will figure out all the other things by itself.

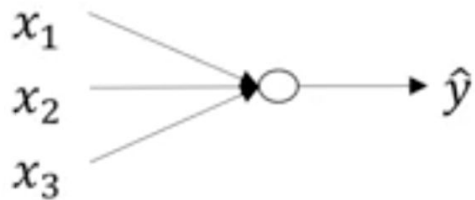
Fully connected or densely connected

This is what we will actually implement

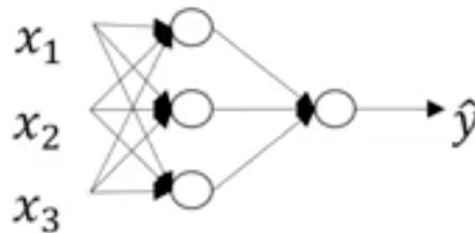


Instead of we deciding that the last hidden unit is “cleanliness”, which depends only on postal code and wealth, we leave it on the neural network to decide what the last hidden unit would be. We just provide it with all input features.

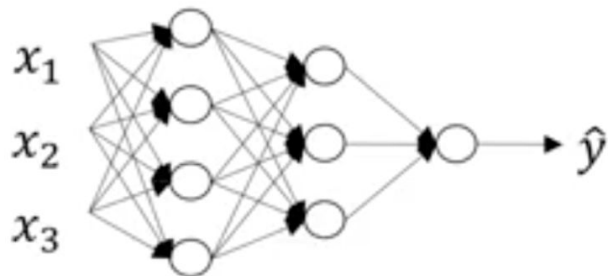
neural network: shallow vs deep



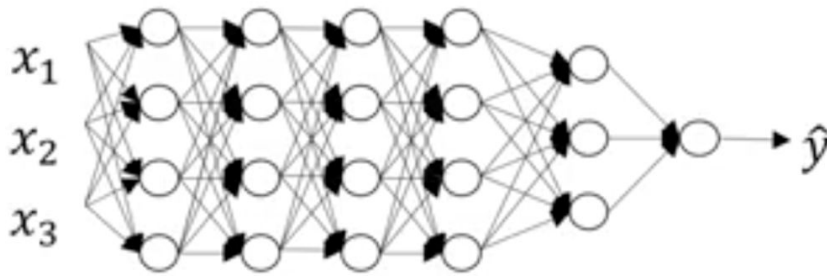
logistic regression



1 hidden layer



2 hidden layers



5 hidden layers