# Introduction to machine learning II

Associated notebook: https://github.com/keeeto/reading-ml-chemistry/blob/master/01_classification_decision_tree.ipynb

# Overview Today(ish)

- Decision trees
- Optimisation
- Evaluation
- Overfitting
- Cross-validation
- Ensemble models

Science and
Technology
Facilities Council

# Setting up a notebook

- You will need a Google account to do this
- Go to https://colab.research.google.com/
- Search for https://github.com/keeeto/reading-ml-chemistry

Science and
Technology
Facilities Council

# A working definition

ML = Representation + Evaluation + Optimisation
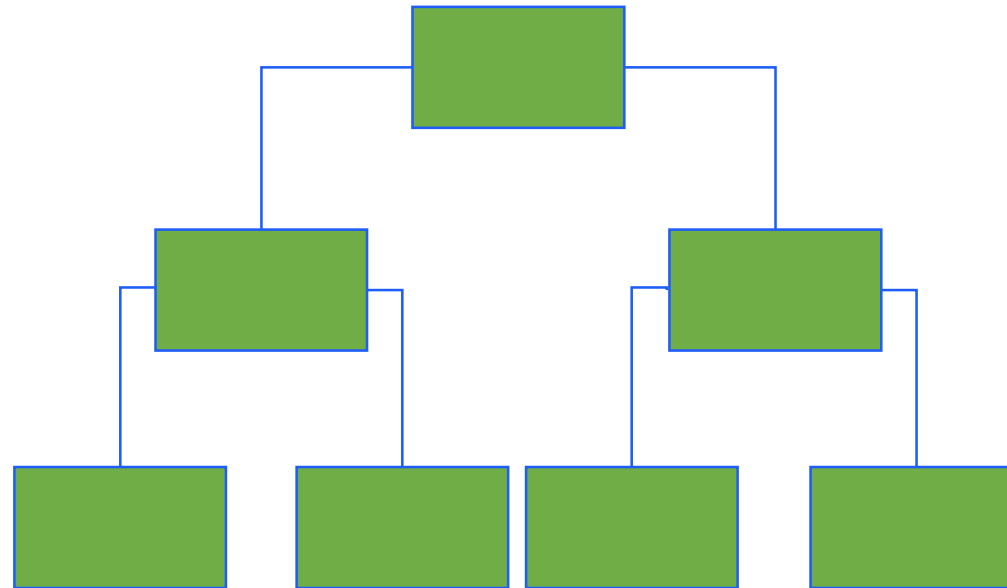
Tapping into the "folk knowledge" needed to advance machine learning applications.

BY PEDRO DOMINGOS

## A Few Useful Things to Know About Machine Learning

"A Few Useful Things to Know About Machine Learning" by Pedro Domingos

# Decision trees



Data is split by features. E.g. brightness of a pixel
Splits are arranged such that the data splits as evenly as possible at each point.

# Decision trees

$$Q_{left}(\theta) = (x, y)|x_f \leq t_j$$
$$Q_{right}(\theta) = Q \setminus Q_{left}(\theta)$$

Data is split according to a threshold value tj.

$$C(Q, \theta) = \frac{n_{left}}{N_j}H(Q_{left}(\theta)) + \frac{n_{right}}{N_j}H(Q_{right}(\theta))$$

The cost of the split is calculated based on some impurity function H() e.g. RMSD of the data.

$$\theta^* = \underset{\theta}{\text{argmin}}\, C(Q, \theta)$$

The splitting parameters are chosen to minimise C at each split.

Science and
Technology
Facilities Council
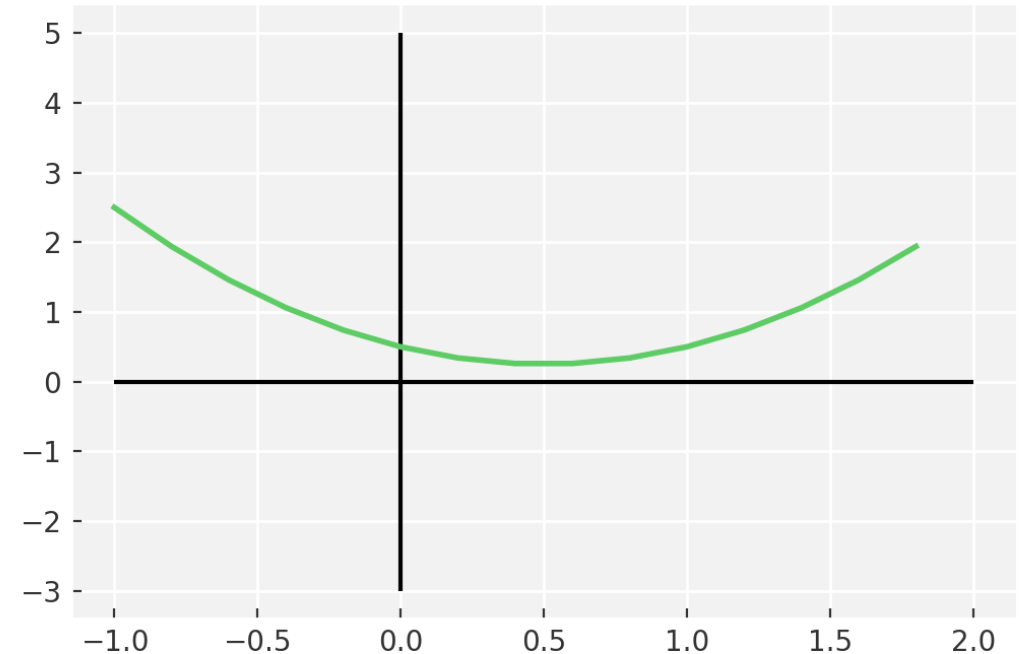
# Go to notebook

# Optimisation/Evaluation

- Evaluation
  - Objective function or scoring function.
  - Distinguish good from bad models.

- Objective function = loss function = cost function
  - Must faithfully represent the "goodness" of a model in a single number

# Evaluation metrics

- Mean squared error
- Used in regression
- Square endures a single minimum
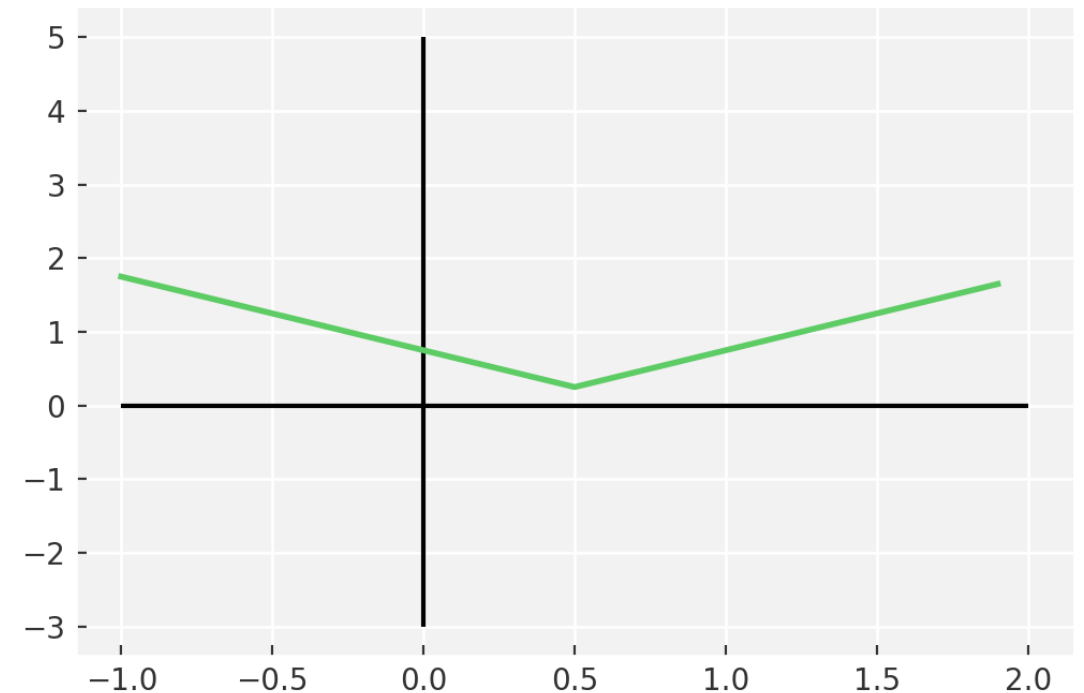- Avoids local minima trapping
- Easy to calculate

$$MSE = \frac{1}{N} \sum (f_i - y_i)^2$$

# Evaluation

- Mean Absolute Error

  - Similar to MSE

  - No quadric term

  - More robust to outliers

  - MSE penalises large differences much more than MAE

  - Large gradients close to zero - slow to optimise
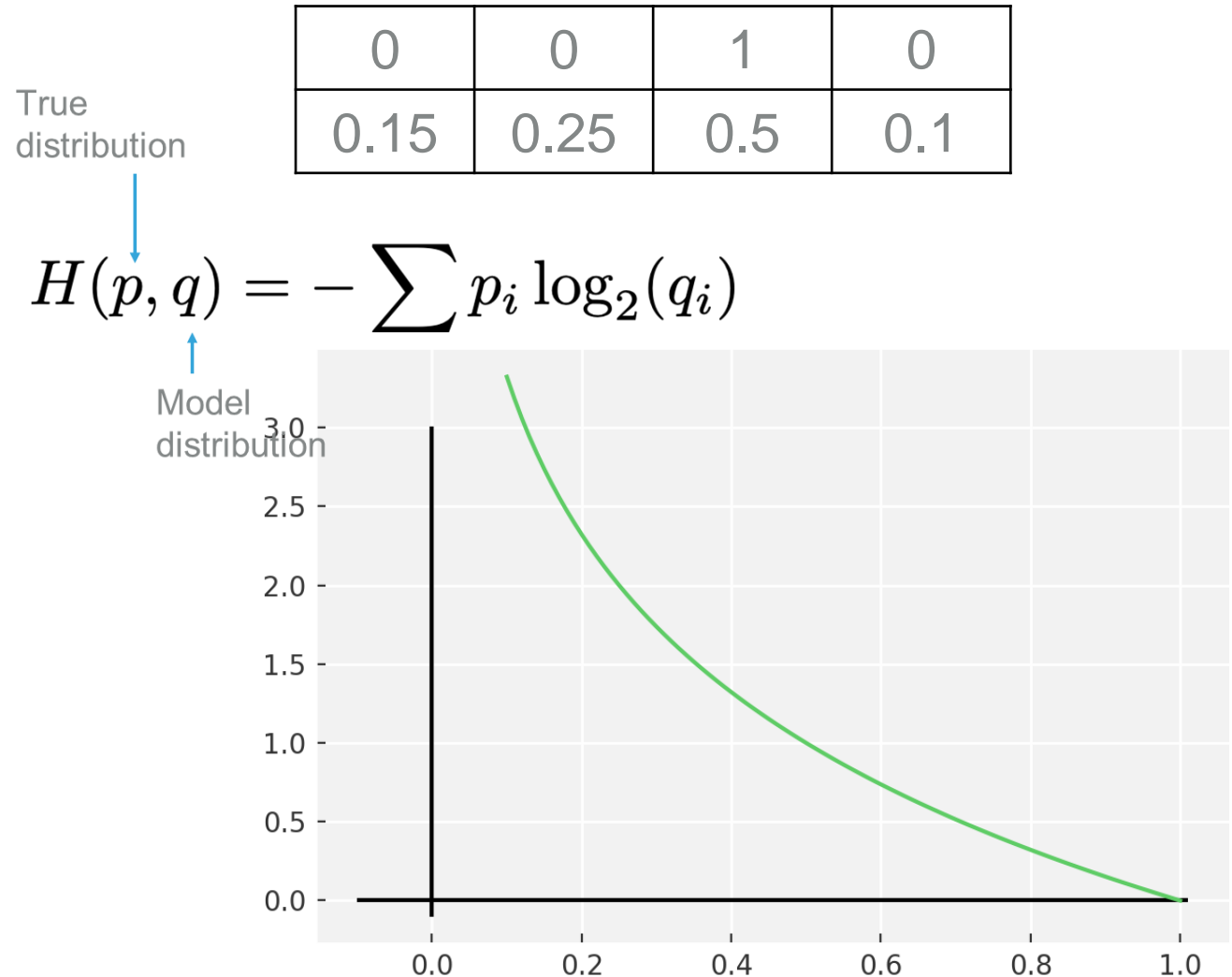
$$MAE = \frac{1}{N} \sum |f_i - y_i|$$

# Evaluation

- Huber loss
- Quadratic close to the minimum
- Linear far from the minimum
- Overcomes problems of MSE and MAE
- More expensive to calculate

$$L_\delta(y, f(x)) = \begin{cases} \frac{1}{2}(y - f(x))^2 & \text{for} |y - f(x)| \leq \delta, \\ \delta |y - f(x)| - \frac{1}{2}\delta^2 & \text{otherwise.} \end{cases}$$



UK RI — Science and Technology Facilities Council

# Evaluation

- Cross entropy
- Used for classification problems
- Tells us how similar our model distribution is to the true distribution
- Penalises all errors, but especially those that are most inaccurate

| 0 | 0 | 1 | 0 |
|------|------|-----|-----|
| 0.15 | 0.25 | 0.5 | 0.1 |

True distribution

$$H(p, q) = -\sum p_i \log_2(q_i)$$

Model distribution

# Evaluation

- Hinge loss
- Used for classification
- Does not seek to reproduce the distribution of data
- 0 as long as the classification is correct

Label(+/-1)
↓

$$L = max(0, 1 - t \cdot y)$$

↑
Prediction

Science and
Technology
Facilities Council

# Evaluation: Table of confusion



https://en.wikipedia.org/wiki/Confusion_matrix
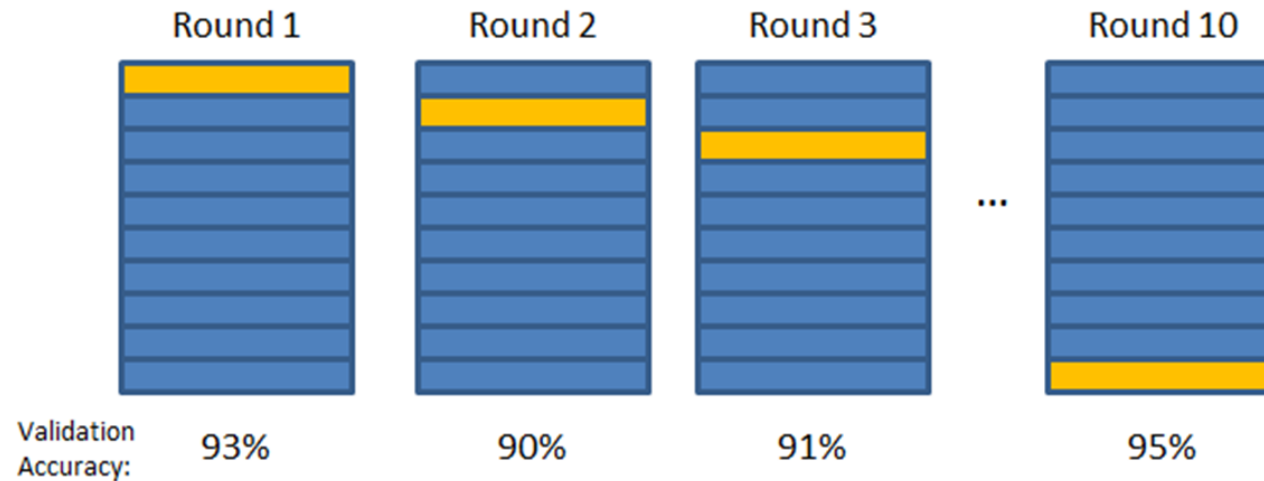
# Evaluation

- Over/under fitting

# Test and validation sets

- The model must always be validated on data not used for testing

- Often something like 20% of data is used for validation

- Make sure that validation and training distributions are the same

# Evaluation

- n-fold cross validation
  - Ensure training/test splits

# Building block Cross Validation

```
from sklearn.model_selection import
cross_val_score
clf = svm.SVC(kernel='linear', C=1)
scores = cross_val_score(clf, X, y, cv=5)
```
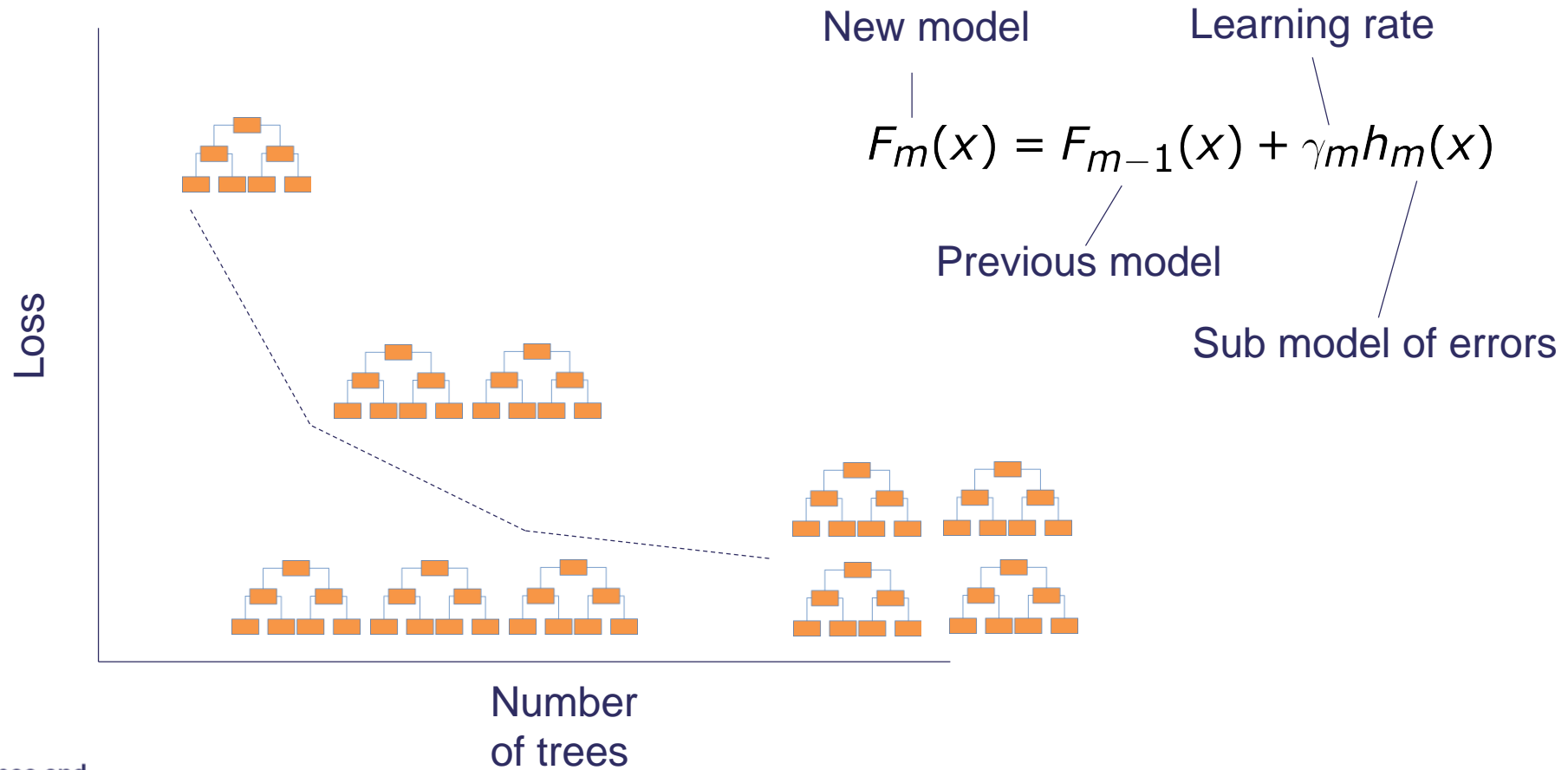
# Go To Notebook

# Boosting + Bagging

- To overcome the limitations of a weak learner we can use booting or bagging.
- Both methods use an ensemble of weak learners to build a strong learner
- Boosting – choose next learner based on the errors of the last learner (gradient boosted decision trees)
- Bagging – stochastically choose next learners (random forests)

# Boosted Decision Trees



New model

Learning rate

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x)$$

Previous model

Sub model of errors

Loss

Number of trees

# Building block: Boosted decision tree

```python
from sklearn import ensemble

gbr = ensemble.GradientBoostingRegressor(loss='lad', max_depth
= 10, learning_rate = 0.015, min_samples_split = 50,
min_samples_leaf = 1, max_features = len(cols), subsample =
0.9, n_estimators = 300)

gbr.fit(X, y)
```

# Go To Notebook

# Concept checklist

- Supervised/unsupervised machine learning
- Classical machine learning/deep learning
- Parameters/hyperparameters
- Features and feature engineering
- Decision trees
- Overfitting
- Evaluation/metrics
- Test/train split, cross-validation
- Bagging and boosting

# Thank you