# Foundational Machine Learning
## Week 3 : Random Forests and Bias Variance Tradeoff

Rohit Babbar
rb2608@bath.ac.uk



UNIVERSITY OF
BATH

- Last week:
  - Decision trees
    - Input = real
    - Output = categorical (discrete)
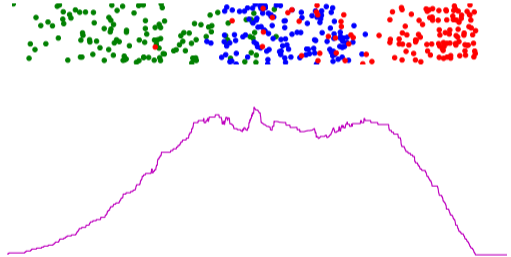
- Last week:
  - Decision trees
    - Input = real
    - Output = categorical (discrete)

- This week:
  - Decision trees with some further input and output types
  - Bias-variance tradeoff – some theoretical insights
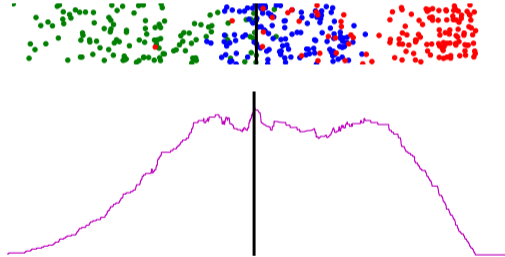  - Bagging - an ensembling technique
  - **Random forests**

Inputs

- Choose best Gini impurity / info gain for all axes

  (1 axis shown; vertical offset for visualisation only)

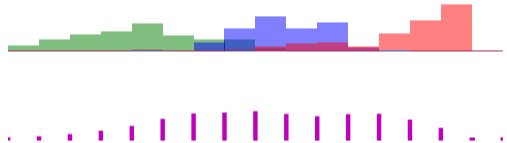- Choose best Gini impurity / info gain for all axes

  (1 axis shown; vertical offset for visualisation only)

- Continuous data may be **quantised**
- e.g. "What is your age?"
  - From 12–17, 18–24 etc.
  - Or year only
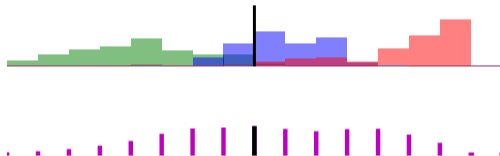
- Continuous data may be **quantised**
- e.g. "What is your age?"
  - From 12–17, 18–24 etc.
  - Or year only

- Split between bins

  (Information gain shown as spikes
  as only defined at bin transitions)

- Similar to quantised (histogram of catgories). . .
  . . . but unordered
- e.g. "What is your favourite cheese?"
- Spliting no longer makes sense!

- Similar to quantised (histogram of catgories). . .
    . . . but unordered
- e.g. "What is your favourite cheese?"
- Spliting no longer makes sense!

- Two choices :
    - Try every assignment of category to the left/right side and pick best
    (one category always goes left to account for symmetry)
    - One category goes down left branch, rest go right

# Categorical input

- Similar to quantised (histogram of catgories)...
  ...but unordered
- e.g. "What is your favourite cheese?"
- Spliting no longer makes sense!

- Two choices :
  - Try every assignment of category to the left/right side and pick best
    (one category always goes left to account for symmetry)
  - One category goes down left branch, rest go right

- One left, rest right is preferred:
  - Fixed storage
  - Simpler code
  - Combinations to test grows linearly
    ($\mathcal{O}(n)$, not $\mathcal{O}(2^{n-1})$; where $n$ = numbers of catgories)

Outputs

- Classification:
  - Split to minimise Gini impurity or maximise information gain
  - Leaf gives answer as most common class to reach it

- Classification:
  - Split to minimise Gini impurity or maximise information gain
  - Leaf gives answer as most common class to reach it

- Regression:
  - Split to minimise variance or maximise information gain
  - Leaf gives answer as mean value to reach it
    (median may confer an advantage — any idea when?)

- Otherwise identical!

# Variance reduction

- Variance of output measures how consistent a node is. . .
    . . . so choose splits that minimise it

- Variance of output measures how consistent a node is. . .
     . . . so choose splits that minimise it

- Variance of left node:
$$\sigma_l^2 = \mathbb{E}\left[(Y_l - \mathbb{E}\left[Y_l\right])^2\right]$$

  similarly for right, $\sigma_r^2$ ($Y_l$ = data that goes left)

- Variance of output measures how consistent a node is...
    ...so choose splits that minimise it

- Variance of left node:
$$\sigma_l^2 = \mathbb{E}\left[(Y_l - \mathbb{E}[Y_l])^2\right]$$

similarly for right, $\sigma_r^2$ ($Y_l$ = data that goes left)

- Minimise weighted combination:
$$L(\texttt{split}) = \frac{n_l}{n}\sigma_l^2 + \frac{n_r}{n}\sigma_r^2$$

$n$ = total exemplar count
$n_l$ = exemplars traveling down left branch
$n_r$ = exemplars traveling down right branch

- Information gain is not often used for regression tasks with decision trees, as it is not immediately applicable

---

[1]https://www.biopsychology.org/norwich/isp/chap8.pdf

- Information gain is not often used for regression tasks with decision trees, as it is not immediately applicable

In order to apply,

- First fit Gaussian distribution to output variable
- Then, compute entropy [1]

$$\frac{1}{2} \log \left(2\pi e \sigma^2\right)$$

- Information gain, thus, is

$$I(\texttt{split}) = \frac{1}{2} \log \left(2\pi e \sigma_p^2\right) - \frac{n_l}{2n} \log \left(2\pi e \sigma_l^2\right) - \frac{n_r}{2n} \log \left(2\pi e \sigma_r^2\right)$$

(same variables as previous slide, with $p$ subscript for parent)

---

[1] https://www.biopsychology.org/norwich/isp/chap8.pdf

From Decision Trees to Random Forests via Bias-variance tradeoff

# Bias-Variance Tradeoff
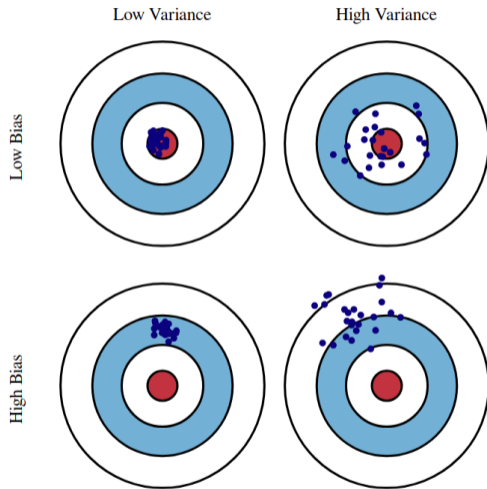


Figure: Pictorial depiction of the components of bias-variance tradeoff

### Definition
An estimator $\hat{\theta}$ of a population parameter $\theta$ is **unbiased** if:

$$\mathbb{E}[\hat{\theta}] = \theta$$

That is, on average, it neither overestimates nor underestimates $\theta$.

Let $X_1, X_2, \ldots, X_n$ be i.i.d. random variables with:

$$\mathbb{E}[X_i] = \mu \quad \text{and} \quad Var(X_i) = \sigma^2$$

The sample mean is defined as:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^{n} X_i$$

Using the linearity of expectation:

$$\mathbb{E}[\bar{X}] = \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n} X_i\right] = \frac{1}{n}\sum_{i=1}^{n}\mathbb{E}[X_i] = \frac{1}{n} \cdot n\mu = \mu$$

Conclusion

$$\boxed{\mathbb{E}[\bar{X}] = \mu}$$

The sample mean is therefore an **unbiased estimator** of the population mean.

Biased estimator (uses $n$ in denominator)

$$\tilde{s}^2 = \frac{1}{n} \sum_{i=1}^{n} (X_i - \bar{X})^2$$

**Key identity:**

$$\sum_{i=1}^{n} (X_i - \bar{X})^2 = \sum_{i=1}^{n} (X_i - \mu)^2 - n(\bar{X} - \mu)^2$$

**Take expectations (i.i.d., $E[X_i] = \mu$, $\mathrm{Var}(X_i) = \sigma^2$):**

$$\mathbb{E}\left[\sum_{i=1}^{n} (X_i - \mu)^2\right] = n\sigma^2, \qquad \mathbb{E}\left[n(\bar{X} - \mu)^2\right] = n \cdot \mathrm{Var}(\bar{X}) = n \cdot \frac{\sigma^2}{n} = \sigma^2$$

**Therefore:**

$$\mathbb{E}\left[\tilde{s}^2\right] = \frac{1}{n}\left(n\sigma^2 - \sigma^2\right) = \left(1 - \frac{1}{n}\right)\sigma^2 = \frac{n-1}{n}\,\sigma^2$$

Turns out that using $(n-1)$ in the denominator makes it unbiased

Consider regression problem with squared error

- Given a training set $D = \{(x_i, y_i)\}_{i=1}^{n}$ such that $x_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$.

- Assume that $y_i = f(x_i) + \epsilon$, where $\epsilon$ is a random variable representing noise with mean 0 and variance $\sigma^2$

- Since $f(.)$ is unknown, we try to approximate it using the training data, and let $\hat{f}(.)$ denote our approximation. Imagine $\hat{f}(.)$ to be decision tree (for regression) that we constructed,

- In regression, this is done by minimising the squared error between $y$ and $f(x)$, i.e. $(y - \hat{f}(x))^2$, where the pair $(x, y)$ could be a training data point or a novel/unseen test data point.

---

[2]Based on Wikipedia article `https://en.wikipedia.org/wiki/Bias-variance_tradeoff`
[3]Derivation as a whole is non-examinable, but you should still have an understanding of the individual parts and the overall idea of the proof

Irrespective of how the classifier $\hat{f}(.)$ is learnt on the data, it's **expected** error on an unseen sample (test) sample $x$ can be decomposed as follows :

$$\mathbb{E}_{D,\varepsilon}\left[\left(y - \hat{f}(x;D)\right)^2\right] = \left(\text{Bias}_D\left[\hat{f}(x;D)\right]\right)^2 + \text{Var}_D\left[\hat{f}(x;D)\right] + \sigma^2$$

where

$$\text{Bias}_D\left[\hat{f}(x;D)\right] \triangleq \mathbb{E}_D\left[\hat{f}(x;D)\right] - f(x)$$

$$\text{Var}_D\left[\hat{f}(x;D)\right] \triangleq \mathbb{E}_D\left[\left(\mathbb{E}_D[\hat{f}(x;D)] - \hat{f}(x;D)\right)^2\right]$$

$$\sigma^2 = \mathsf{E}_y\left[\left(y - f(x)\right)^2\right]$$

Note $\triangleq$ means this is a definition

- Rewriting the LHS of the equation on the previous slide as MSE (Mean squared error)

$$\text{MSE} \triangleq \mathbb{E}\Big[\big(y - \hat{f}(x)\big)^2\Big] \qquad \textit{writing the expectation } \mathbb{E} \textit{ without subsripts}$$

$$= \mathbb{E}\Big[\big(f(x) + \varepsilon - \hat{f}(x)\big)^2\Big] \qquad \text{since } y \triangleq f(x) + \varepsilon$$

$$= \mathbb{E}\Big[\big(f(x) - \hat{f}(x)\big)^2\Big] + 2\,\mathbb{E}\Big[\big(f(x) - \hat{f}(x)\big)\varepsilon\Big] + \mathbb{E}[\varepsilon^2]$$

- Using the fact that for independent r.v. $X$ and $Y$, $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$

$$\mathbb{E}\Big[\big(f(x) - \hat{f}(x)\big)\varepsilon\Big] = \mathbb{E}\big[f(x) - \hat{f}(x)\big]\,\mathbb{E}\big[\varepsilon\big] \quad \text{since } \varepsilon \text{ is independent from } x$$

$$= 0 \qquad\qquad\qquad \text{since } \mathbb{E}\big[\varepsilon\big] = 0$$

- Expanding the first term below :

$$\mathbb{E}\Big[\big(f(x) - \hat{f}(x)\big)^2\Big] = \mathbb{E}\Big[\big(f(x) - \mathbb{E}[\hat{f}(x)] + \mathbb{E}[\hat{f}(x)] - \hat{f}(x)\big)^2\Big]$$

$$= \mathbb{E}\Big[\big(f(x) - \mathbb{E}[\hat{f}(x)]\big)^2\Big] + 2\,\mathbb{E}\Big[\big(f(x) - \mathbb{E}[\hat{f}(x)]\big)\big(\mathbb{E}[\hat{f}(x)] - \hat{f}(x)\big)\Big]$$

$$+ \mathbb{E}\Big[\big(\mathbb{E}[\hat{f}(x)] - \hat{f}(x)\big)^2\Big]$$

$$\mathbb{E}\Big[\big(f(x) - \mathbb{E}[\hat{f}(x)]\big)^2\Big] = \mathbb{E}\big[f(x)^2\big] - 2\,\mathbb{E}\Big[f(x)\,\mathbb{E}[\hat{f}(x)]\Big] + \mathbb{E}\Big[\mathbb{E}[\hat{f}(x)]^2\Big]$$

$$= f(x)^2 - 2\,f(x)\,\mathbb{E}[\hat{f}(x)] + \mathbb{E}[\hat{f}(x)]^2$$

$$= \Big(f(x) - \mathbb{E}[\hat{f}(x)]\Big)^2$$

The term in red :

$$\mathbb{E}\Big[\big(f(x) - \mathbb{E}[\hat{f}(x)]\big)\big(\mathbb{E}[\hat{f}(x)] - \hat{f}(x)\big)\Big] = \mathbb{E}\Big[f(x)\,\mathbb{E}[\hat{f}(x)] - f(x)\hat{f}(x) - \mathbb{E}[\hat{f}(x)]^2 + \mathbb{E}[\hat{f}(x)]\,\hat{f}(x)\Big]$$

$$= f(x)\,\mathbb{E}[\hat{f}(x)] - f(x)\,\mathbb{E}[\hat{f}(x)] - \mathbb{E}[\hat{f}(x)]^2 + \mathbb{E}[\hat{f}(x)]^2$$

$$= 0$$

$$\mathsf{MSE} = \Big(f(x) - \mathbb{E}[\hat{f}(x)]\Big)^2 + \mathbb{E}\Big[\big(\mathbb{E}[\hat{f}(x)] - \hat{f}(x)\big)^2\Big] + \sigma^2$$

$$= \mathsf{Bias}\,\big(\hat{f}(x)\big)^2 + \mathsf{Var}\,\big[\hat{f}(x)\big] + \sigma^2$$

# Bias-variance tradeoff - V

- Variance: Captures how much the classifier changes if trained on a (slightly) different training set. How "over-specialized" is it to a particular training set (overfitting)?

- Bias: What is the inherent error that classifiers incurs even with infinite training data? This is due to the classifier being "biased" to a particular kind of solution (e.g. linear classifier).

- Noise: How much is the data-intrinsic noise? It's a measure of the ambiguity due to the data distribution and feature representation ($\epsilon$ in the above example).

- The goal in Random Forest (next) is to lower the variance $\mathbb{E}\left[\left(\mathbb{E}\left[\hat{f}(x)\right] - \hat{f}(x)\right)^2\right]$ by having more and more estimates of $\hat{f}(x)$, i.e. growing lots of decision trees into (random) forests !

Random Forests

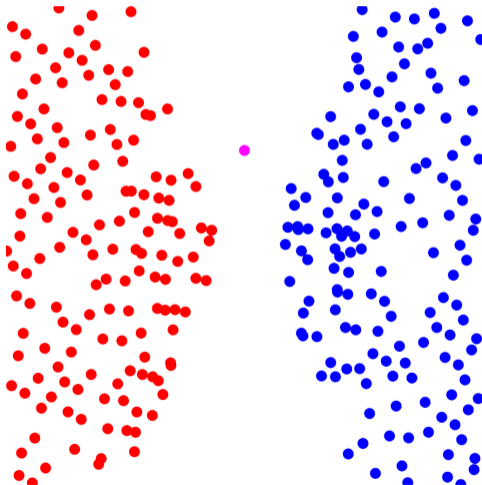- Random forest = decision trees + ensemble learning

- Random forest = decision trees + ensemble learning

- Ensemble learning = combining multiple estimators
    - Different models (e.g. linear regression, decision tree, SVM, neural network)
                    or
    - Same model, randomised training so each is different

    (using estimator to distinguish from model)

- Random forest = decision trees + ensemble learning

- Ensemble learning = combining multiple estimators
  - Different models (e.g. linear regression, decision tree, SVM, neural network)
    or
  - Same model, randomised training so each is different

  (using estimator to distinguish from model)

- Random forest:
  - Many decision trees (hence name)
  - Randomised training using **bagging**,
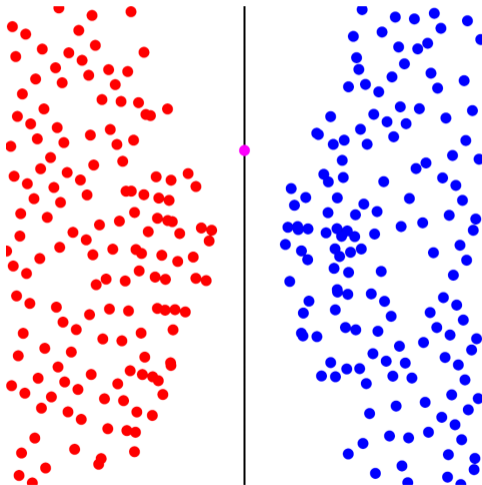    a specific ensemble learning technique
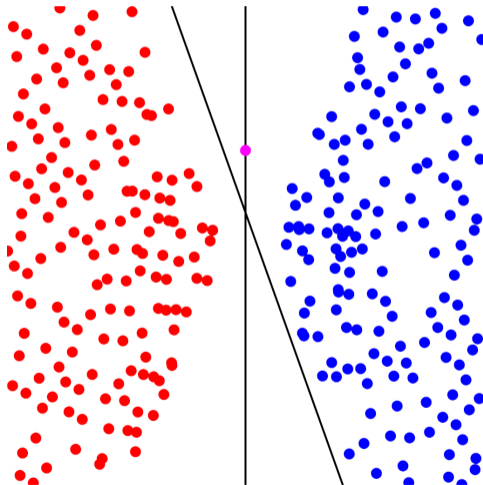
- Which class (red or blue) should the magneta dot be?

- Which class (red or blue) should the magneta dot be?
- An obvious classification boundary ...
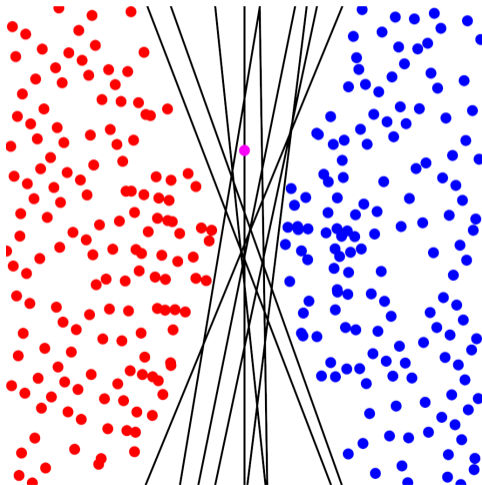
- Which class (red or blue) should the magneta dot be?
- An obvious classification boundary . . .
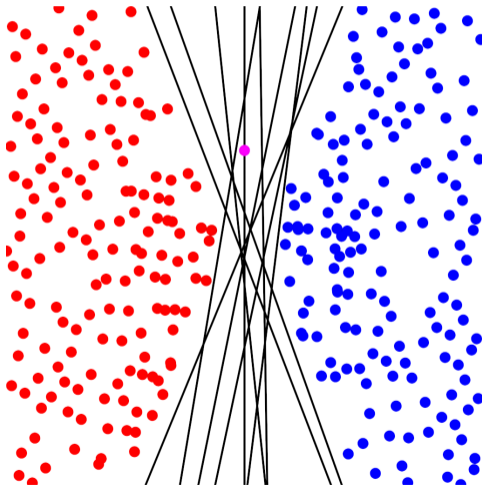- But this is just as good (suggesting blue)

- Which class (red or blue) should the magneta dot be?
- An obvious classification boundary . . .
- But this is just as good (suggesting blue)
- As are all of these!

- Which class (red or blue) should the magneta dot be?
- An obvious classification boundary ...
- But this is just as good (suggesting blue)
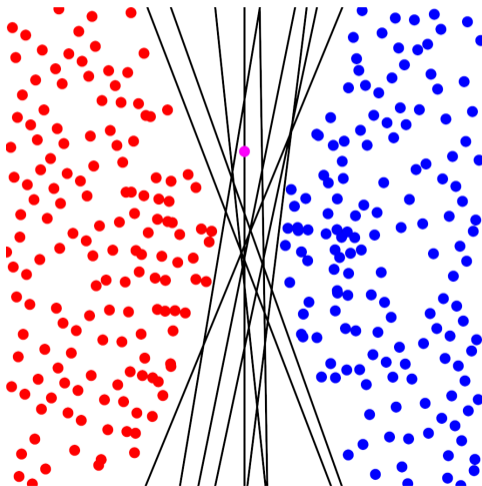- As are all of these!

- Models can be fit in many ways due to:
  - Insufficient data
  - Noisy data
  - **Model not complex enough**
    (curved boundaries can also separate this data!)

- Which class (red or blue) should the magneta dot be?
- An obvious classification boundary . . .
- But this is just as good (suggesting blue)
- As are all of these!

- Models can be fit in many ways due to:
  - Insufficient data
  - Noisy data
  - **Model not complex enough**
    (curved boundaries can also separate this data!)

- Ensembles have many estimators. . .
  . . . to capture this ambiguity

- Ensembles need **diversity**

- Estimators must make **different mistakes**
  i.e. if all make same mistake $\implies$ ensemble will repeat it
- Increasing estimator diversity at expense of individual performance $\rightarrow$ better ensemble!
  (up to a limit)

- However, constructing an ensemble would require more data
- Instead of collecting more data. . .
  . . . fake it from available data

- However, constructing an ensemble would require more data
- Instead of collecting more data...
    ...fake it from available data

- Given data set of size $n$:
    Create new data set by drawing, with replacement, $n$ times
    (there will be repetitions of data-points due to replacement)

- An ensemble technique!

- Short for "Bootstrap AGGregatING"
- Bagging = Bootstrapping applied to estimator output
  (via optimised parameters)

- An ensemble technique!

- Short for "Bootstrap AGGregatING"
- Bagging = Bootstrapping applied to estimator output
  (via optimised parameters)

- Algorithm:
  1. Select $S$, size of ensemble
  2. Create $S$ bootstrap draws of original data set
  3. Train estimator on each
  4. Combine outputs of all estimators for each query

- Another ensemble technique!

- Bootstrap applied to features
    i.e. fit each estimator with a random subset of features

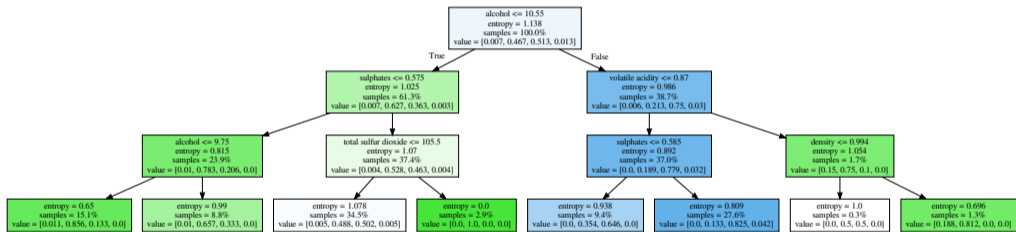- For decision tree: new bootstrap for each split

- Random forest = decision trees + bagging + random subspace method

- Algorithm:
  1. Select $S$, number of trees (more is better, up to a limit)
  2. Create $S$ bootstrap draws of original data set
  3. Train decision tree on each, with random subspace method
  4. Combine outputs of all trees for each query
     - Classification : The decision trees vote – ensemble outputs winner
     - Regression : Take the mean/median of all of the estimates
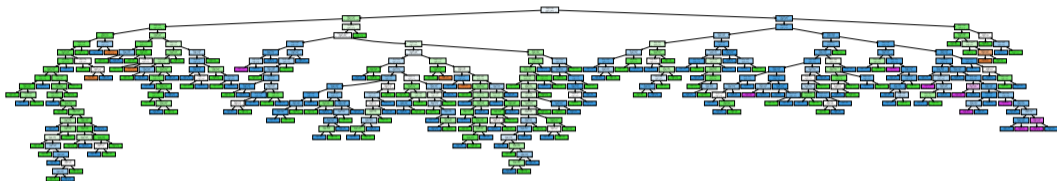
Explainability

- 1599 exemplars, split: 1199 train, 400 test.
- Input: 11 measurable features:
  - fixed acidity
  - volatile acidity
  - citric acid
  - residual sugar
  - chlorides
  - free sulfur dioxide
  - total sulfur dioxide
  - density
  - pH
  - sulphates
  - alcohol

- Output: 1–10 human rating
    (reduced to 1–4 here, to fit on screen)

- Can be used for classification or regression!
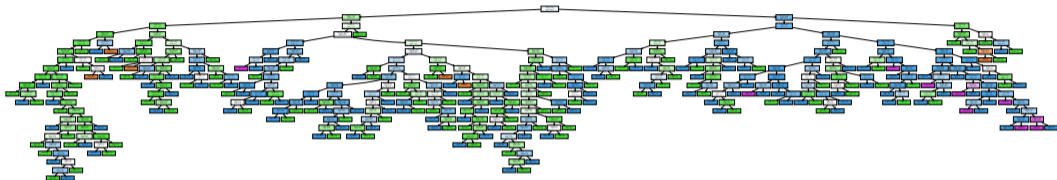
Classification tree, max depth 3:



Accuracy = 69% (and explainable)

- A deeper classification tree:
  - Accuracy = 71% (somewhat less explainable)

- A deeper classification tree:
  - Accuracy = 71% (somewhat less explainable)

- Random forests with 32 trees
  - Classification: Accuracy = 79% (better than 71%)
  - Impossible to visualise/understand!

- "Upgraded" decision trees to
  - Handle more types of inputs & outputs

  - Random forests! (still one of the best)
  - Output probabilities

- Notes:
  - One of the fastest algortihms
  - Many variants, e.g. gradient boosting

- Next week :
    Making sure a ML system is working!

- For more kinds of random forest:
  "Decision Forests for Classification, Regression, Density Estimation, Manifold Learning and Semi-Supervised Learning" by **Criminisi, Shotton and Konukoglu**
  `http://research.microsoft.com/apps/pubs/default.aspx?id=155552`

- For more on ensemble methods:
  "Diversity creation methods: a survey and categorisation"
  by **Brown, Wyatt, Harris and Yao**
  `http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.421.349&rep=rep1&type=pdf`

Cheese board, CC Worm That Turned, Attribution-Share Alike 4.0 International,
`https://commons.wikimedia.org/wiki/File:Welsh_cheese_board.JPG`

Wine data set from `https://archive.ics.uci.edu/ml/datasets/wine+quality`