Tree depth 2



| Bedrooms | Bathrooms | Latitude | Longitude | Price |
|----------|-----------|----------|-----------|-------|
| 1 | 1 | 40.71 | -73.94 | 3055 |
| ... | ... | ... | ... | ... |

# Decomposition of Feature Space

Mittwoch, 30. Juni 2021     11:50

## Feature space:

- 2 features: x1, x2: plane
- 3 features: x1, x2, x3: cube
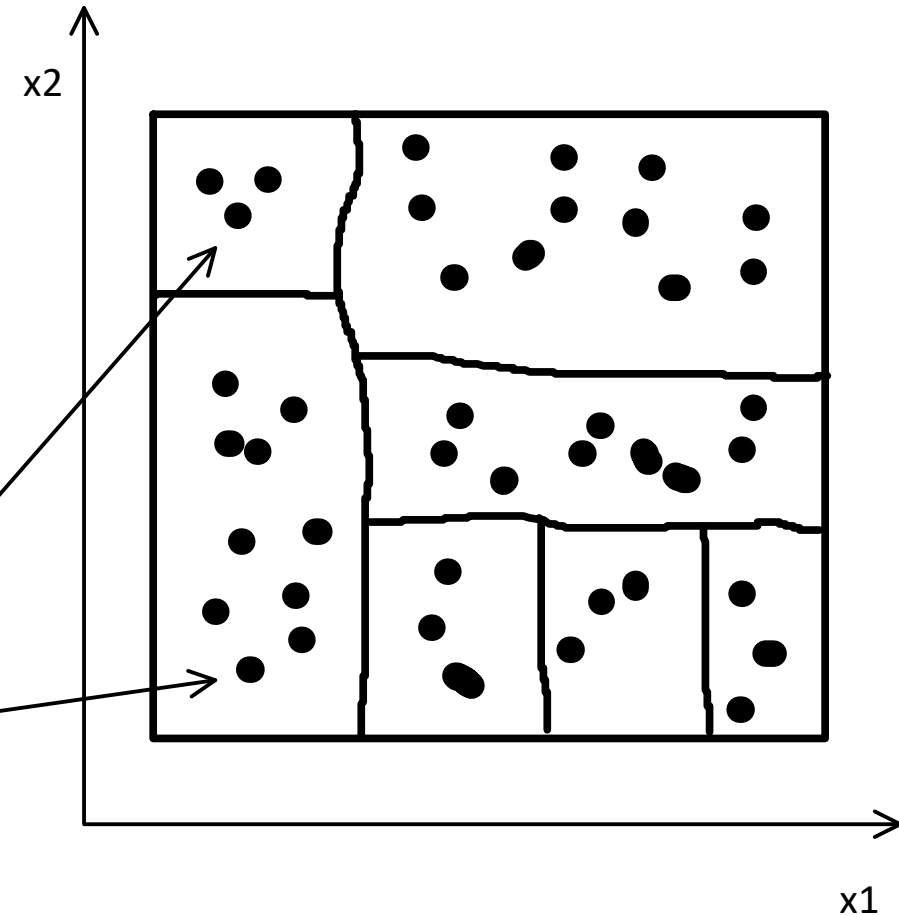- n features: x1, …, xn: n-dimensional cube

## Find rectangular non-overlapping decomposition of space

- Leads to n-dimensional sub-cubes (regions)
- Each region should be as "uniform" as possible
- Average value of all records in region (regression)
- Majority vote in region (classification)

Instances

Finding an optimal decomposition automatically  is infeasable (combinatorical explosion)

## Simplifcation for presentation purposes

- Only two features: x1, x2
- Regions are rectangles

## General idea of decision trees:
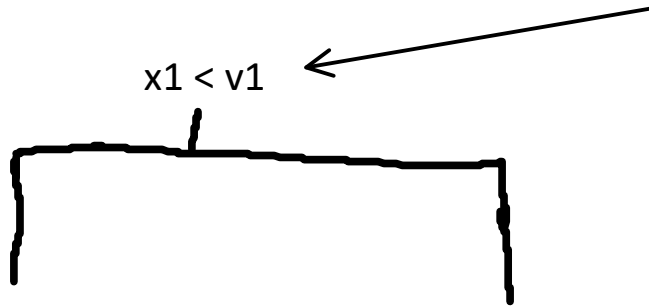
- Gready decomposition of feature space by recursive binary splitting
- Stop splitting according to criterion, e.g.
  - minimal numbers of instances in region
  - max depth of tree
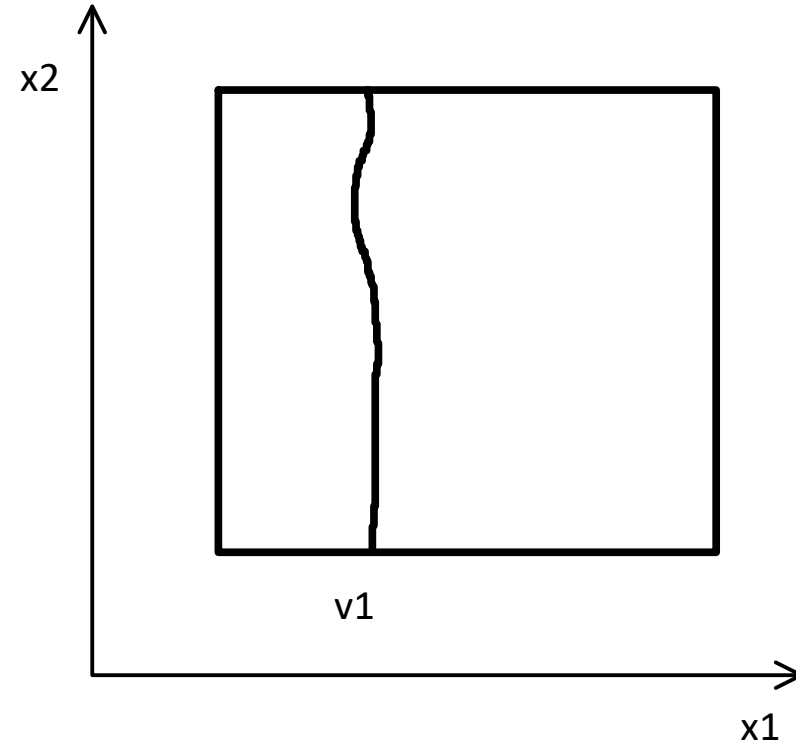  - minimal performance gain

## Finding splits

- Search through all features
- For each feature: consider all split values
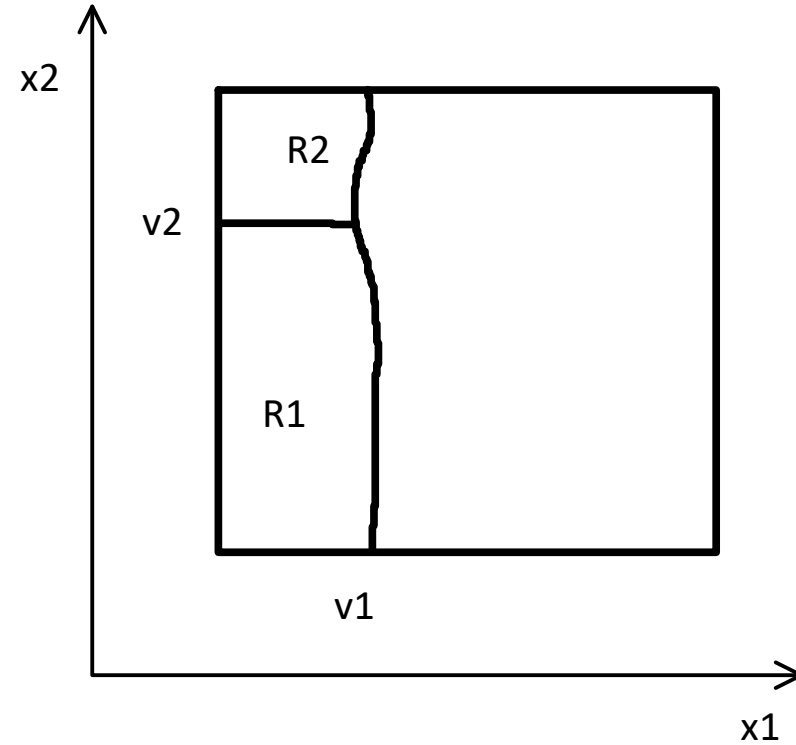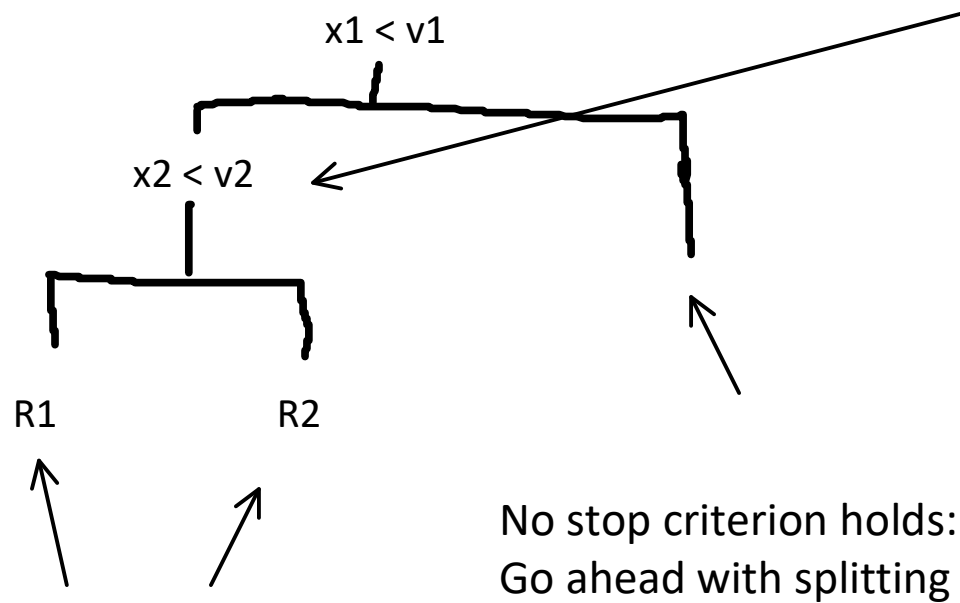- Take "best" (feature, splitval) combination
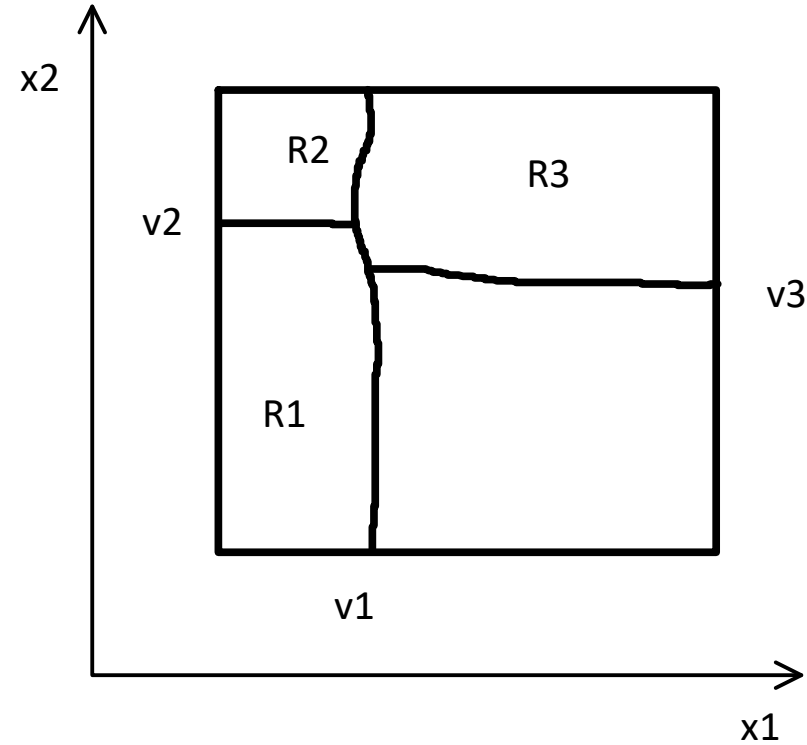
x1 < v1

| split feature | split value |
|---------------|-------------|
| x1            | v1          |

No stop criterion holds:
Go ahead with splitting

x2

v1

x1

| split feature | split value |
|:---:|:---:|
| x2 | v2 |

x1 < v1

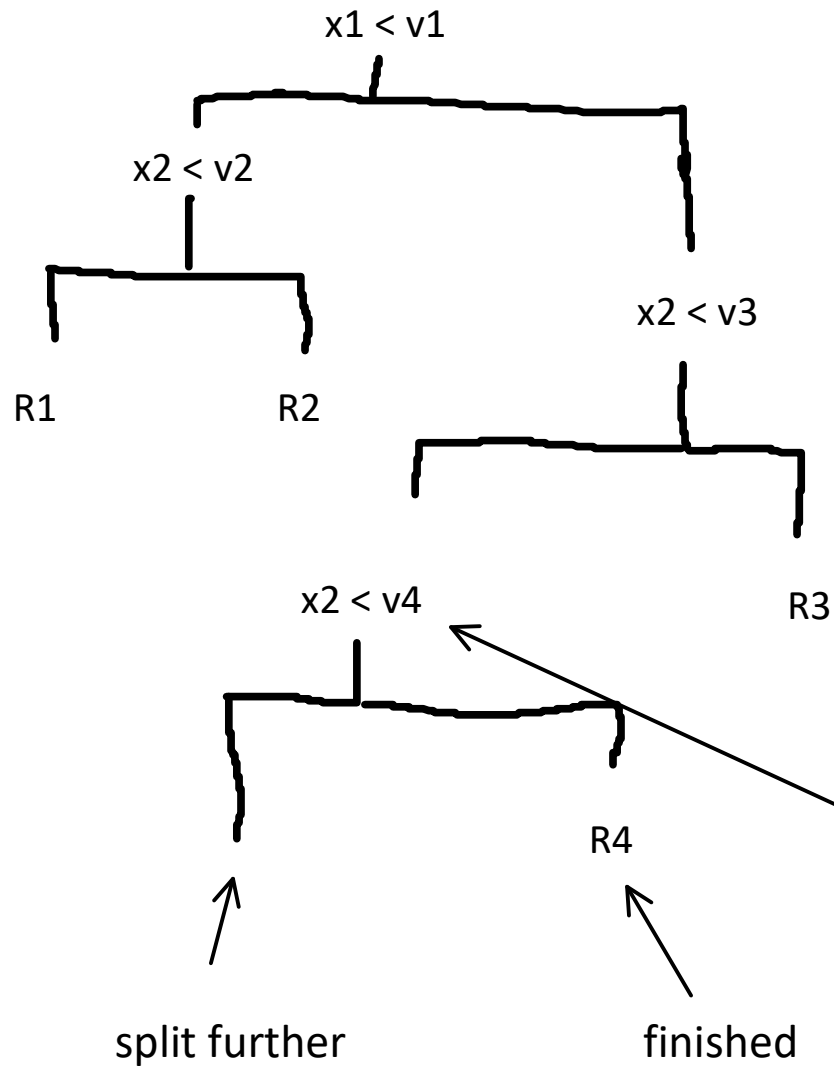x2 < v2

R1          R2

No stop criterion holds:
Go ahead with splitting

Stop criterion holds:
do not further split

x2

v2

R2

R1

v1

x1

| split feature | split value |
|---------------|-------------|
| x2 | v3 |

x1 < v1

x2 < v2

x2 < v3

R1          R2

R3

split further

finished

x2

v2

R2

R3

v3

R1

v1

x1

$x1 < v1$

$x2 < v2$

$x2 < v3$

R1        R2

$x2 < v4$

R3

R4

split further        finished

| split feature | split value |
|---|---|
| x2 | v4 |

| split feature | split value |
|---|---|
| x1 | v5 |

stop splitting

bedrooms
1 | 2, 3, 4

Split on best combination
of feature and split point

...

bathrooms
1, 2 | 3

latitude
< 40.71

latitude
>=40.71

$$\mathrm{MSE}(R) = \frac{1}{\mathrm{count}(y_i)} \sum_{\mathbf{x}_i \in R} (y_i - \mathrm{avg}(y_i))^2$$



|     | x1 | x2 | y  |
|-----|----|----|----|
| **x1** | 1  | 1  | 20 |
| **x2** | 1  | 2  | 30 |
| **x3** | 2  | 3  | 20 |
| **x4** | 3  | 1  | 10 |
| **x5** | 4  | 2  | 20 |

$$\mathrm{avg}(y_i) = 20$$

$$\mathrm{MSE}(R) = \frac{1}{5} \left( (20 - 20)^2 + (30 - 20)^2 + (20 - 20)^2 + (10 - 20)^2 + (20 - 20)^2 \right)$$

$$= \frac{1}{5} (0 + 100 + 0 + 100 + 0)$$

$$= \frac{1}{5} * 200$$

$$= 40$$

$$\text{SplitVariance} = \frac{1}{\text{count}(R_1)}\text{MSE}(R_1) + \frac{1}{\text{count}(R_2)}\text{MSE}(R_2)$$

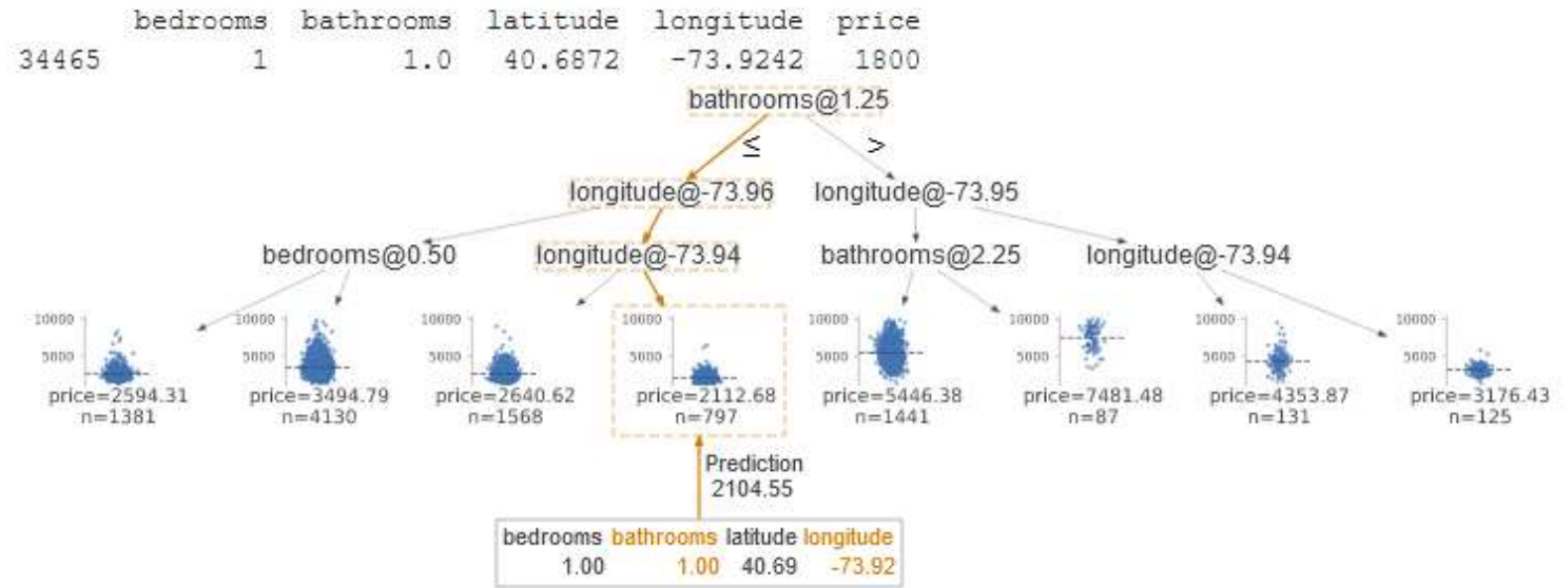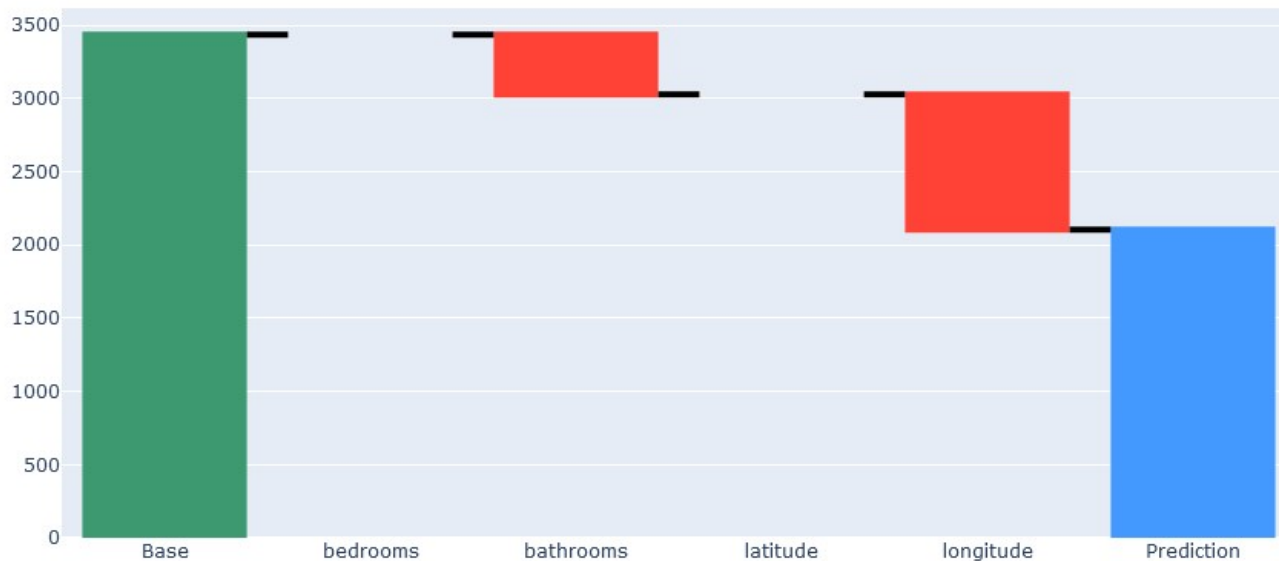Prediction : 2607.8296004300305

Prediction : 2104.545399698341
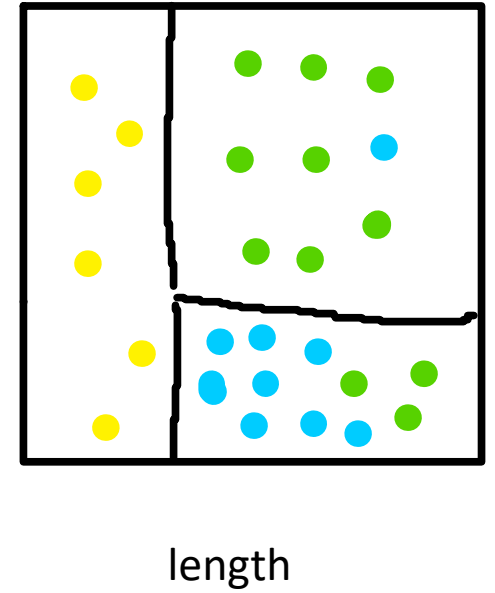
## Tree depth 2



| | petal_length | petal_width | class | target |
|---|---|---|---|---|
| 1 | 1.4 | 0.2 | Iris-setosa | 0 |
| 2 | 1.3 | 0.2 | Iris-setosa | 0 |
| 51 | 4.5 | 1.5 | Iris-versicolor | 1 |
| 52 | 4.9 | 1.5 | Iris-versicolor | 1 |
| 101 | 5.1 | 1.9 | Iris-virginica | 2 |
| 102 | 5.9 | 2.1 | Iris-virginica | 2 |

## Tree depth 3

$$\text{Gini}(R) = 1 - \sum_{c \in C} \text{Prob}(c, R)^2$$

| R | Region |
|------|-------------------|
| Prob | Probability |
| Gini | Gini Value |
| C | Set of all classes |
| c | Single class |

Class tags:
- Class1: orange
- Class2: blue



$$\text{Gini}(R) = 1 - \left( \left( \frac{2}{6} \right)^2 + \left( \frac{4}{6} \right)^2 \right)$$

$$= 1 - \left( \left( \frac{1}{3} \right)^2 + \left( \frac{2}{3} \right)^2 \right)$$

$$= 1 - \left( \frac{1}{9} + \frac{4}{9} \right)$$

$$= 1 - \frac{5}{9}$$

$$= \frac{4}{9}$$

$$\text{InformationGain} = \text{Gini}(R) - \left( \frac{1}{\text{count}(R_1)} \text{Gini}(R_1) + \frac{1}{\text{count}(R_2)} \text{Gini}(R_2) \right)$$

# Example Information Gain

Mittwoch, 30. Juni 2021         18:30

See:

https://towardsdatascience.com/decision-tree-an-algorithm-that-works-like-the-human-brain-8bc0652f1fc6

Prof. Dr. Ingo Claßen

## Problems of trees

- Too specific
- Overfit on trainings data

Assume leafs only contain one record

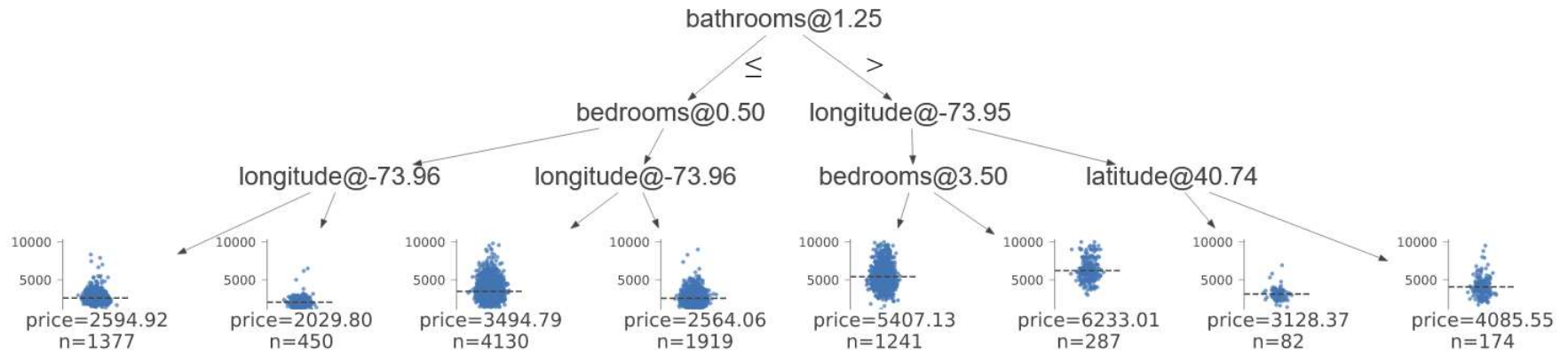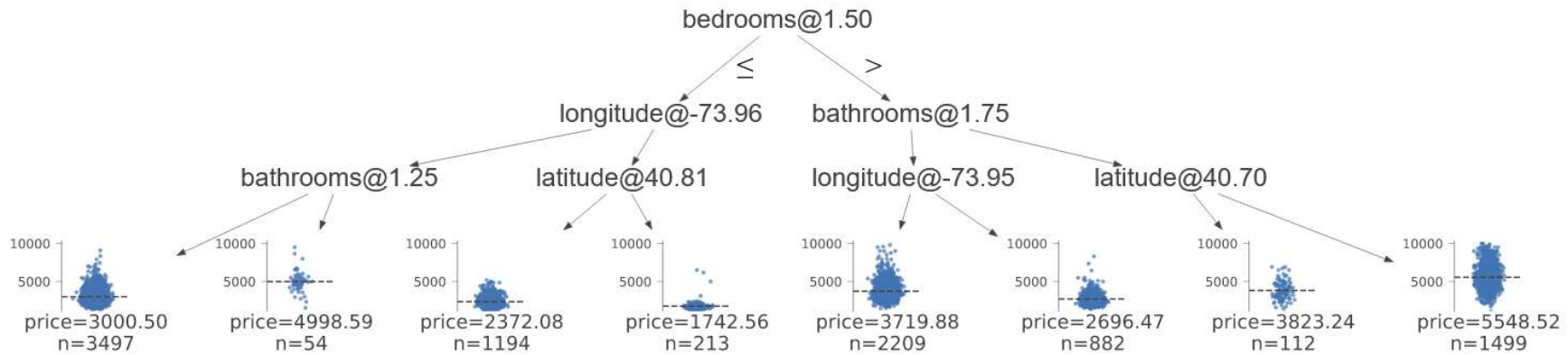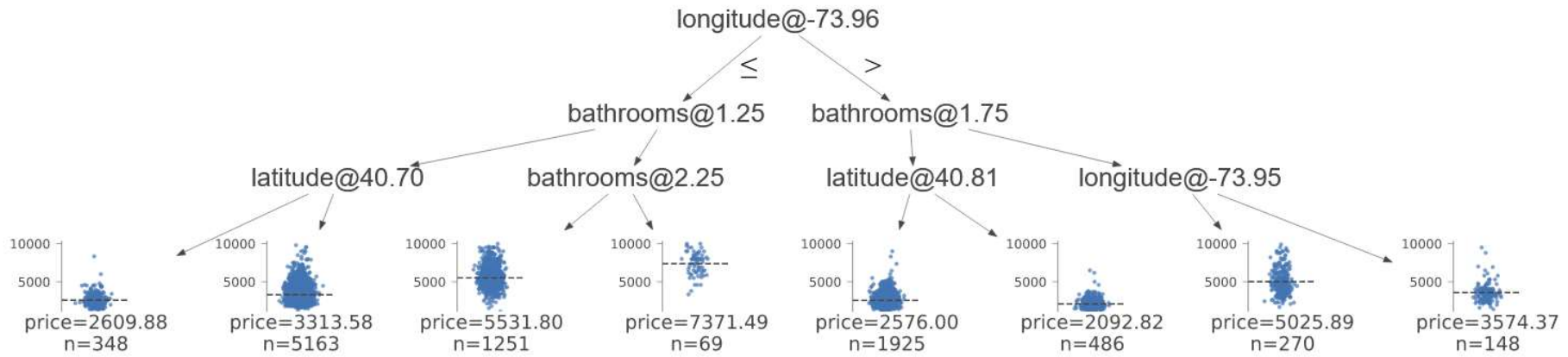- Trainings error would get zero

## Train many tree models

- Randomly select subset of training data (with replacement, bootstrapping)
- Randomly select subset of features

## Prediction - Regression

- Ask all trees for prediction
- Average results
- Two levels of averaging
  - Leaf level of each singular tree
  - Average of all trees

## Prediction - Classification

- Ask all trees for prediction
- Majority vote
- Two levels of voting
  - Leaf level of each singular tree
  - Combining votes of all trees

## Permutation Importance