

# Regression Trees

## Probability and Statistics for Data Science

Carlos Fernandez-Granda



These slides are based on the book [Probability and Statistics for Data Science](#) by Carlos Fernandez-Granda, available for purchase [here](#). A free preprint, videos, code, slides and solutions to exercises are available at <https://www.ps4ds.net>

# Regression

**Goal:** Estimate **response** from **features**

**Optimal estimator:** Conditional mean

**Problem:** Intractable to compute due to curse of dimensionality

# Linear regression

Response  $y$  is approximated as an **linear** (affine) function of the features  $x$

$$y \approx \sum_{i=1}^d \beta[i]x[i] + \alpha$$

**Assumption:** Response increases **or** decreases **proportionally** to each feature (if we fix other features)

# Example

**Response:** Temperature in Manhattan (Kansas)

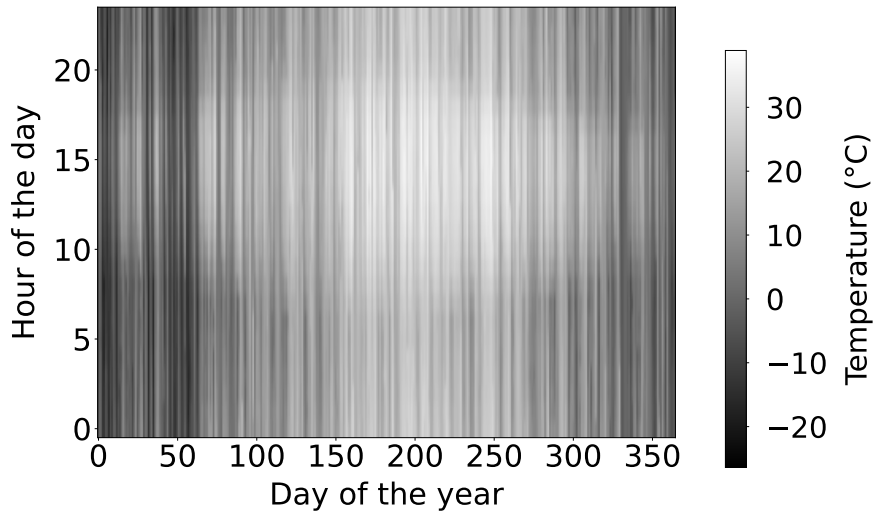
**Features:**

- (1) Hour of the day (0-23)
- (2) Day of the year (1-365)

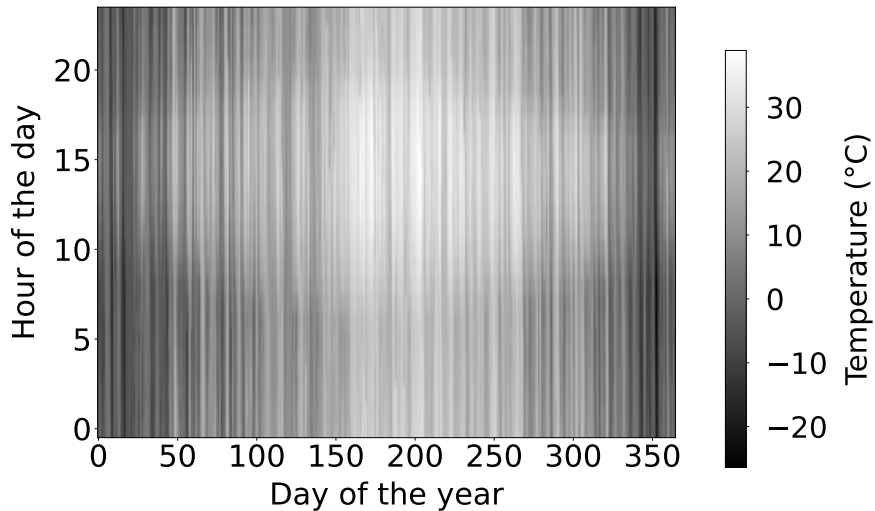
**Training data:** 2015

**Test data:** 2016

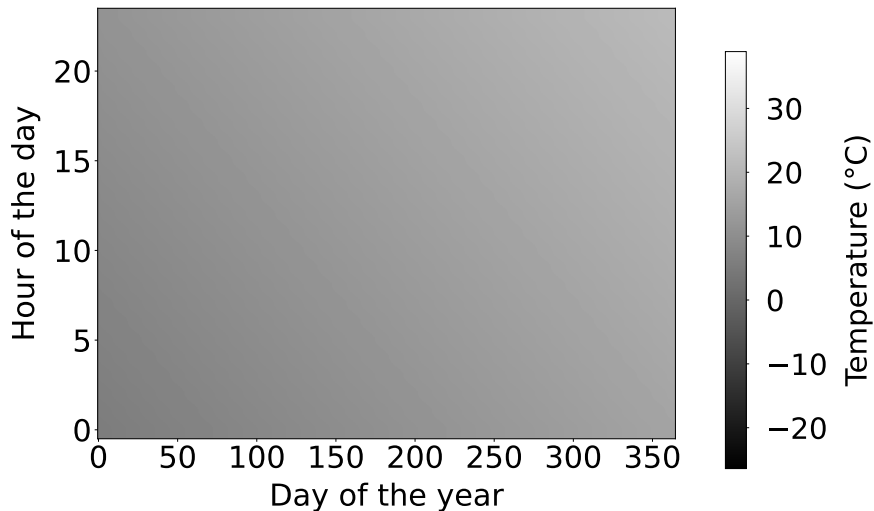
## Training data



## Test data



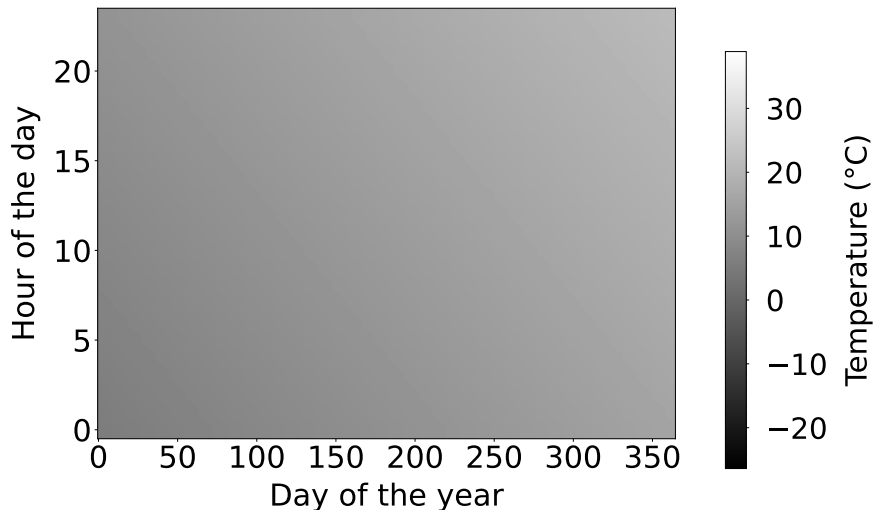
Linear model:  $0.25 \text{ hour} + 0.03 \text{ day} + 5.85$



Response **increases or decreases proportionally** to each feature (if we fix other features)



Linear model:  $0.25 \text{ hour} + 0.03 \text{ day} + 5.85$



Training error:  $10.8^{\circ}\text{C}$

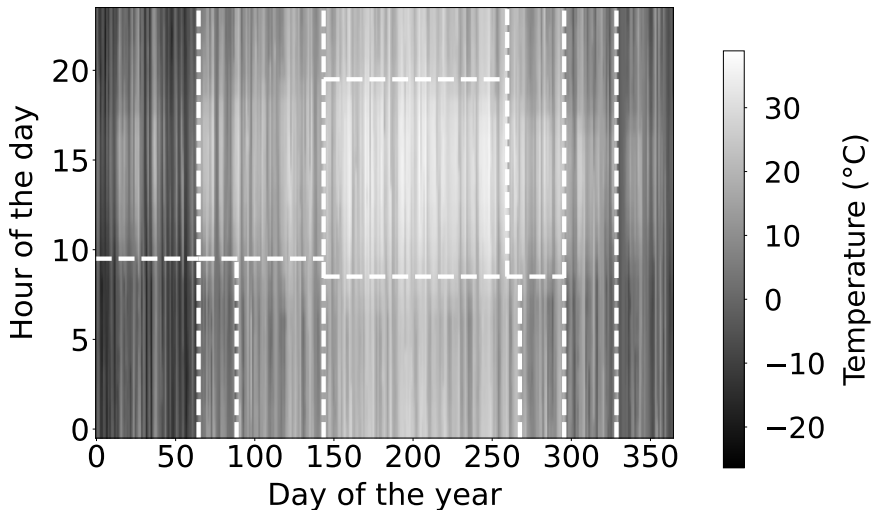
Test error:  $11.0^{\circ}\text{C}$

# Challenge

How to learn **nonlinear** model?

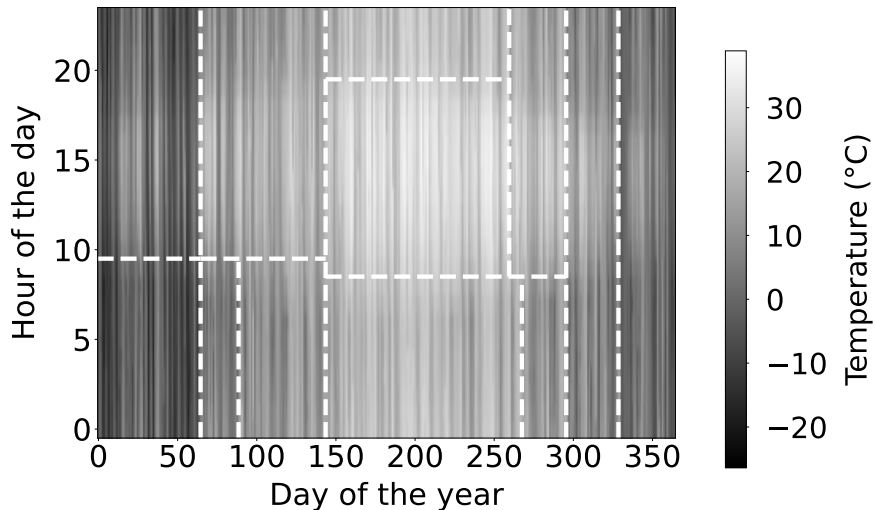
## Idea

(1) Partition feature space into regions



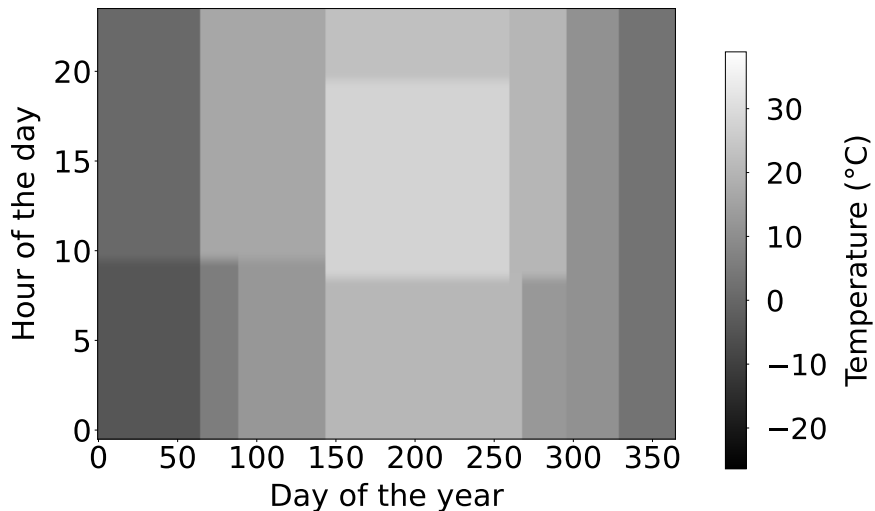
## Idea

(2) Assign constant estimate to each region

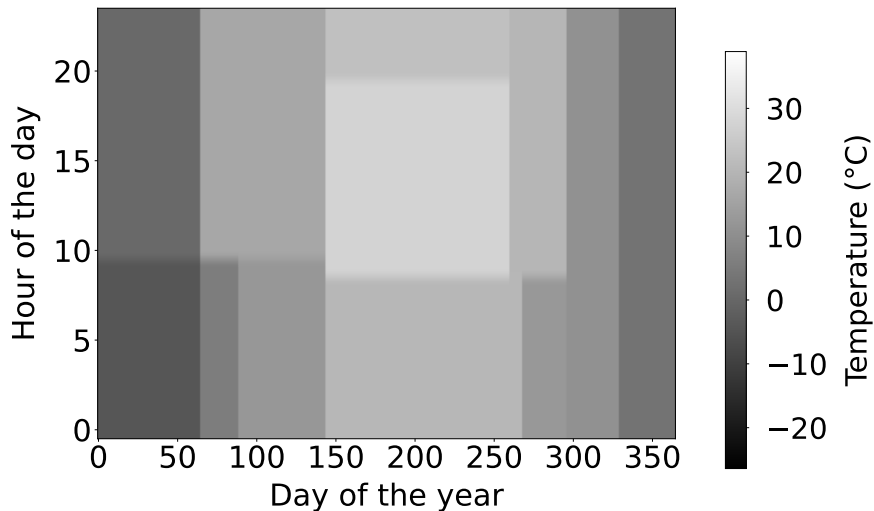


## Idea

(2) Assign constant estimate to each region



Works pretty well!



Training error: 5.5°C

Test error: 6.2°C

## Two key questions

How to compute constant estimate?

How to choose the regions?

## Constant estimate?

Consider the  $n_R$  feature-response pairs  $(x_i, y_i)$  in region  $R$

$$\text{RSS}(\alpha) := \sum_{\{i: x_i \in R\}} (y_i - \alpha)^2$$

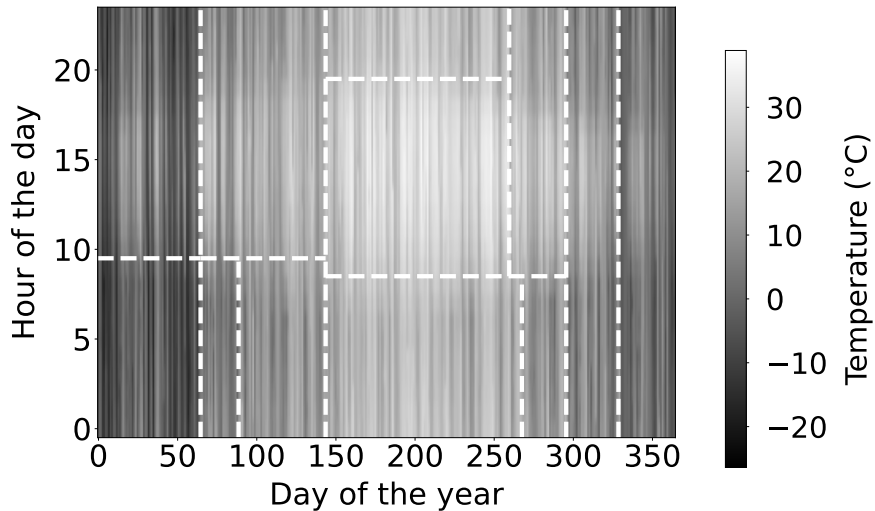
$$\begin{aligned} \frac{d \text{RSS}(\alpha)}{d\alpha} &= - \sum_{\{i: x_i \in R\}} (y_i - \alpha) \\ &= n_R \alpha - \sum_{\{i: x_i \in R\}} y_i \end{aligned}$$

$$\frac{d^2 \text{RSS}(\alpha)}{d\alpha^2} = n_R$$

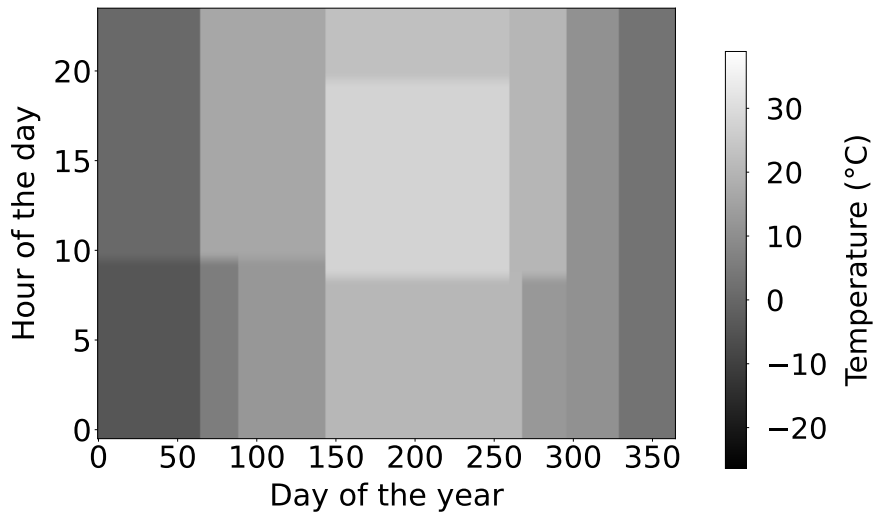
$$\alpha_{\min} = \frac{1}{n_R} \sum_{\{i: x_i \in R\}} y_i$$



Constant estimate?



Just average!

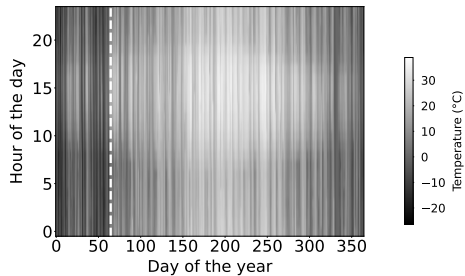
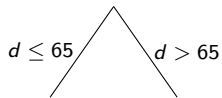


# Regions?

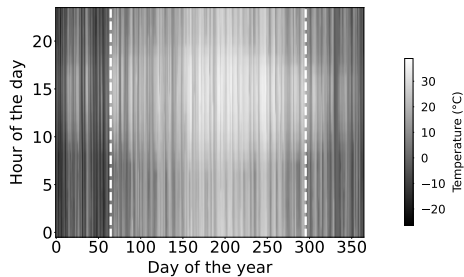
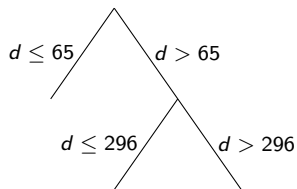
Possible regions explode exponentially with number of features!

**Idea:** Use a binary tree to represent the regions

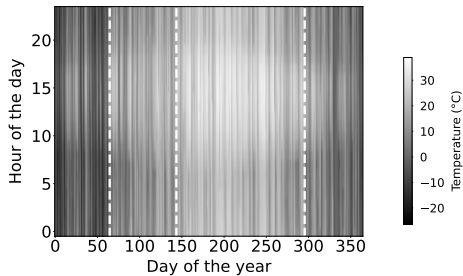
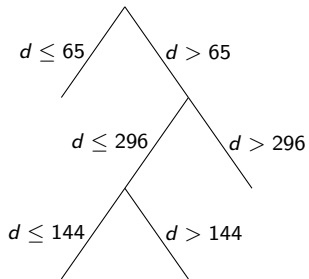
# Tree



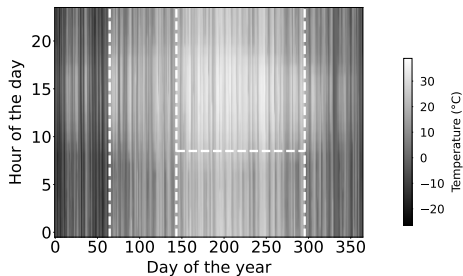
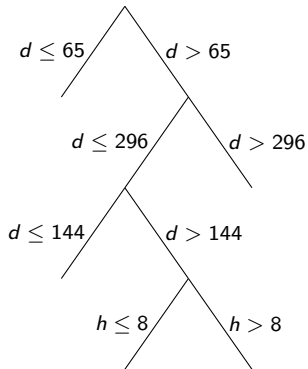
# Tree



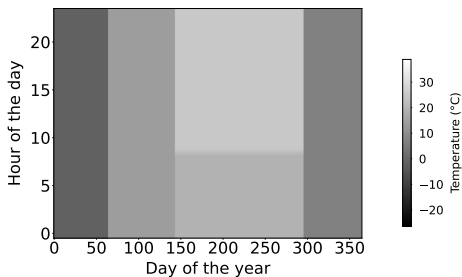
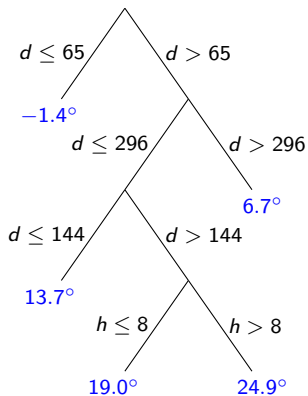
# Tree



# Tree

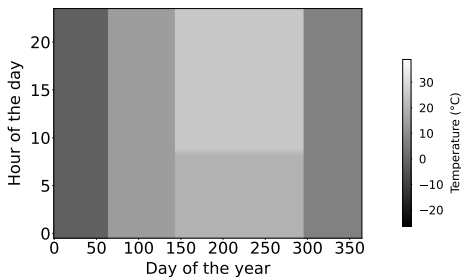
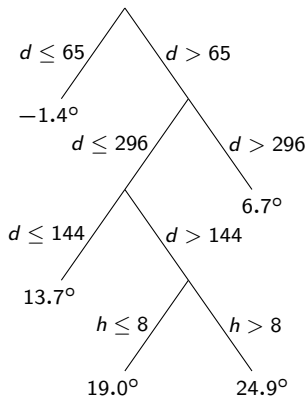


# Region estimates

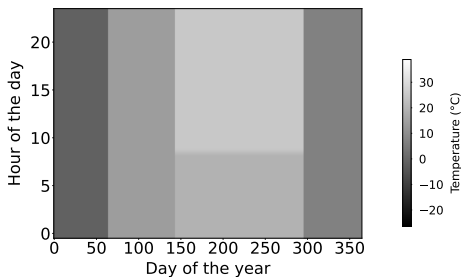
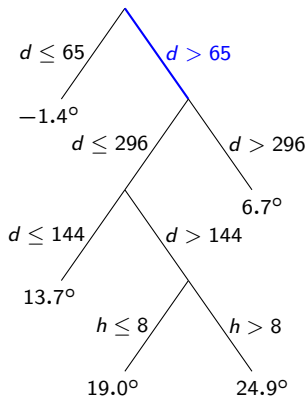




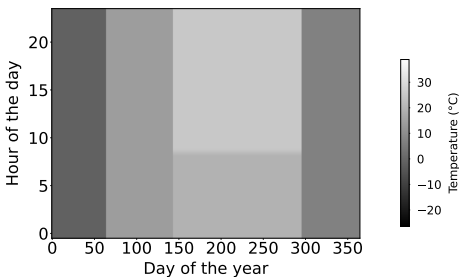
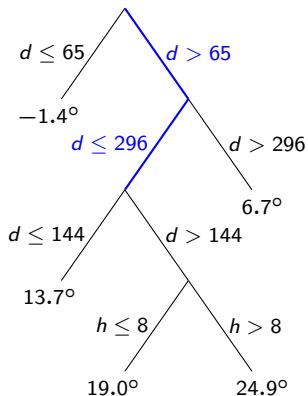
August 19 ( $d := 251$ ) at 3 am ( $h := 3$ )?



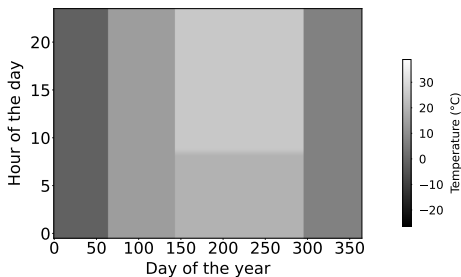
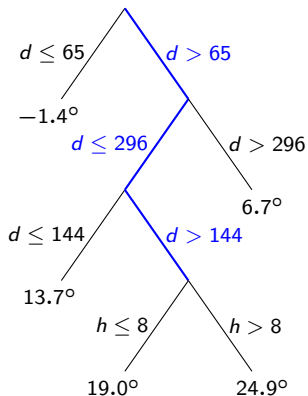
August 19 ( $d := 251$ ) at 3 am ( $h := 3$ )?



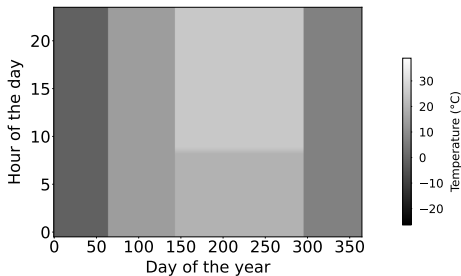
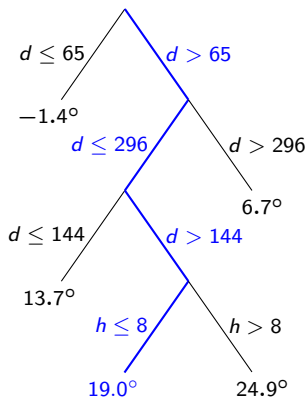
August 19 ( $d := 251$ ) at 3 am ( $h := 3$ )?



August 19 ( $d := 251$ ) at 3 am ( $h := 3$ )?



August 19 ( $d := 251$ ) at 3 am ( $h := 3$ )?



Interpretable!

Actual temperature (in 2023):  $22^\circ$

# How do we build the tree?

**Idea:** Choose tree with smallest training error

Tree with depth  $h$  and  $2^h$  leaves

Number of possible bifurcations?  $b := 2^h - 1$

At each bifurcation (1)  $d$  features and (2)  $t$  thresholds

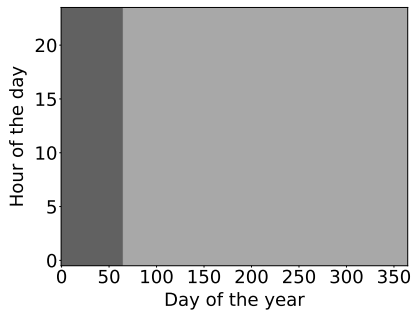
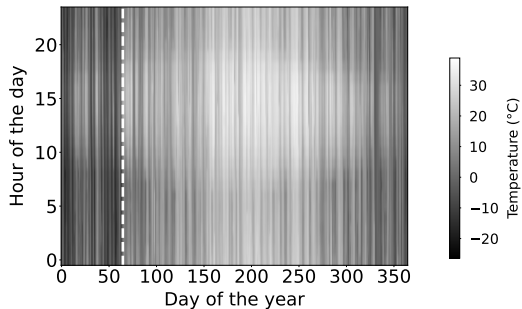
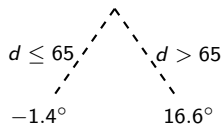
Number of possible trees?  $(dt)^b$

For  $h := 4$ ,  $d := 10$ ,  $t := 100$ :  $10^{45}$  trees!

## Recursive binary splitting

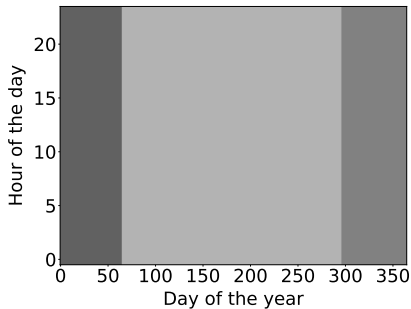
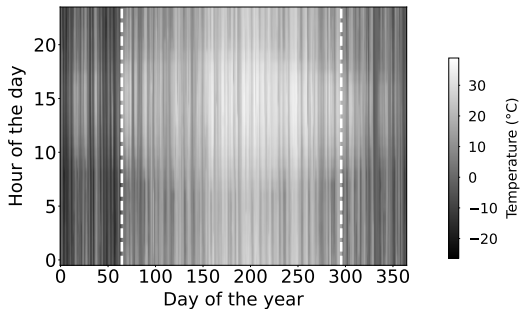
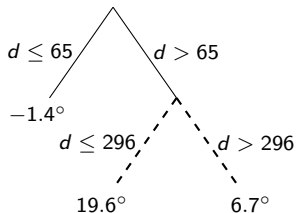
Add one bifurcation at a time, being greedy

# Bifurcation 1

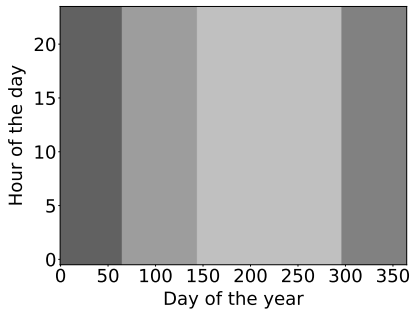
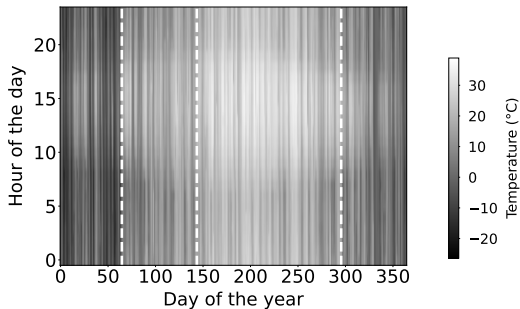
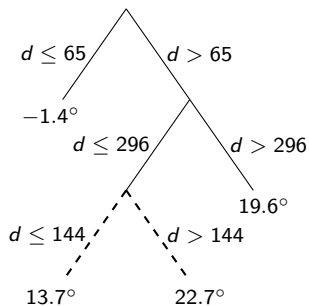




## Bifurcation 2



## Bifurcation 3



# Residual Sum of Squares (RSS)

Data:  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

Regions:  $R_1, \dots, R_m$

Estimates at each region:  $\alpha_1, \dots, \alpha_m$

$$\text{Residual Sum of Squares} := \sum_{r=1}^m \sum_{\{i: x_i \in R_r\}} (y_i - \alpha_r)^2$$

## Choosing a split

$$\text{RSS} := \sum_{r=1}^m \sum_{\{i: x_i \in R_r\}} (y_i - \alpha_r)^2$$

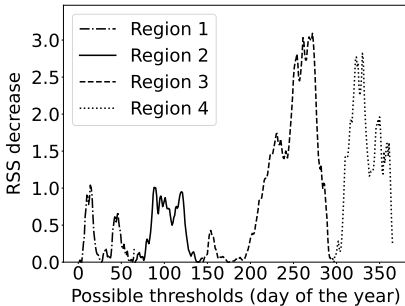
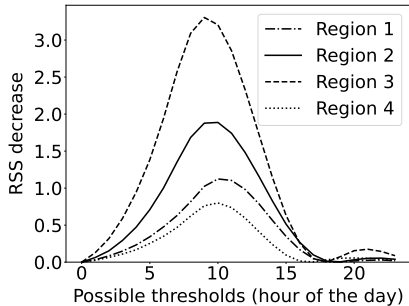
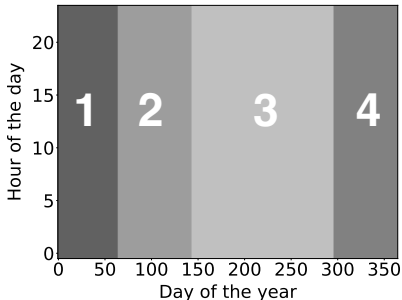
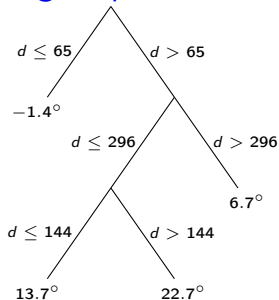
If we split region  $R_r$  into subregions  $A$  and  $B$

$$\begin{aligned} \Delta \text{RSS} := & \sum_{\{i: x_i \in R_r\}}^n (y_i - \alpha_r)^2 - \sum_{\{i: x_i \in A\}}^n (y_i - \alpha_A)^2 \\ & - \sum_{\{i: x_i \in B\}}^n (y_i - \alpha_B)^2 \end{aligned}$$

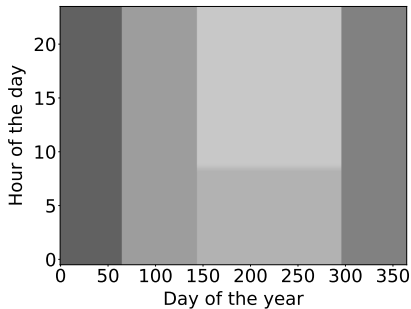
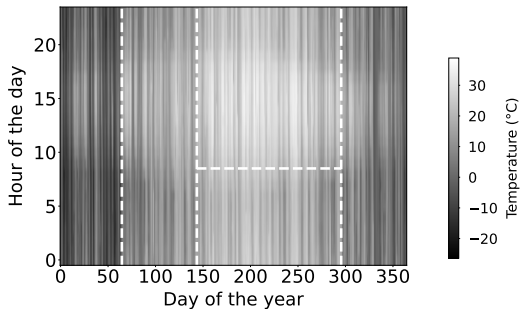
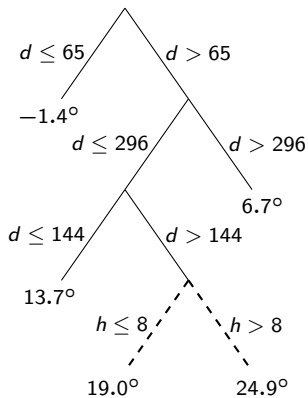
Depends on (1) region, (2) feature and (3) threshold

Choose split that **maximizes**  $\Delta \text{RSS}$

# Choosing a split

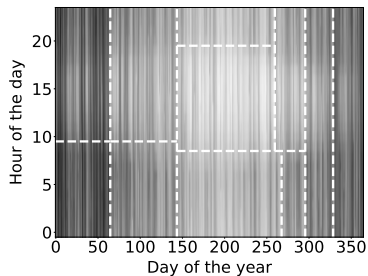


## Bifurcation 4

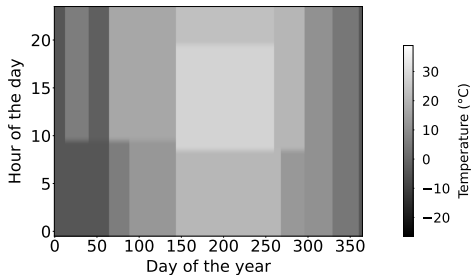
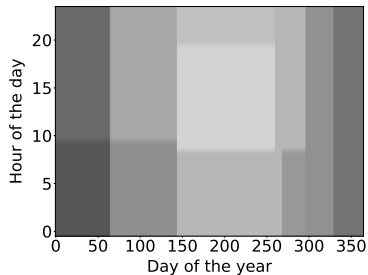
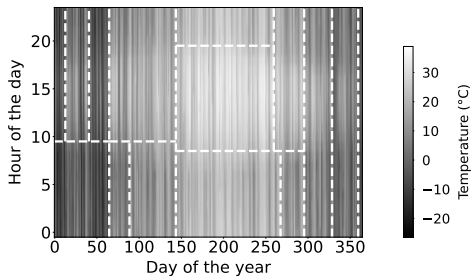


# When to stop?

11 leaves

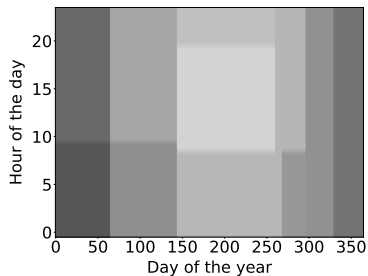


15 leaves

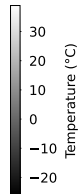
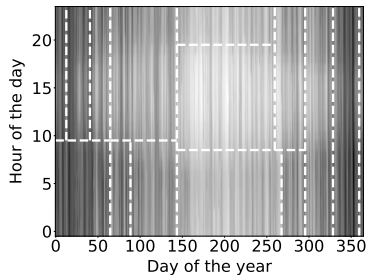
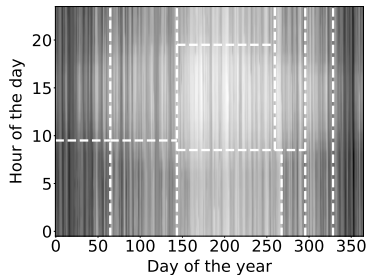
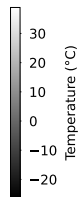
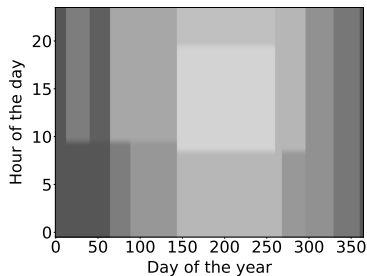


# Test data

## 11 leaves

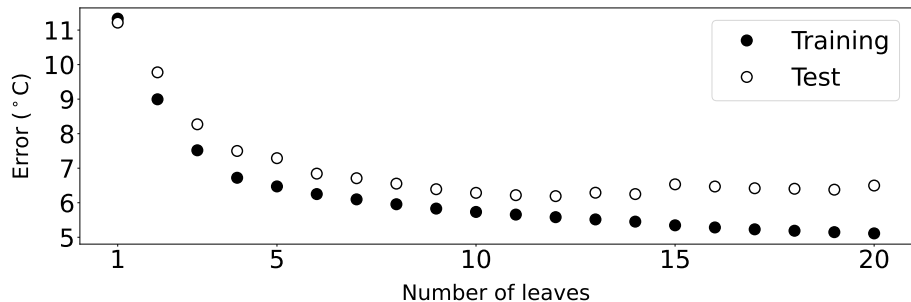


## 15 leaves





## Training and test error



# What have we learned?

How to build nonlinear regression models using trees

To be careful about overfitting!