

DMML, 31 Jan 2019

Supervised learning — Predict from training data

↳ Classification problem

Sometimes — need to predict a numeric value

— Exam marks

— Salary, mileage for a car, ...

Attributes:  $A_i$   $\begin{cases} \text{numeric} \checkmark \\ \text{categorical} \end{cases}$

Mileage of a car: Weight of car, engine size,  
...

Exam marks: Exam marks in practice tests,  
IQ score, height, weight, dob

how to  
rule them out?

Not today

Model  $f(a_1, \dots, a_n) = y$

Variables  $x_1, \dots, x_k$  for  $k$  input attributes

Calculate a function  $f(x_1, \dots, x_k)$  that approximates output well

Assume linear relationship (by inspection?)

$$y = \theta_0 + \theta_1 x_1 + \dots + \theta_k x_k$$

Learning = Parameter Estimation  
- Find  $\theta_i$

Find "best" set of  $\theta_i$

Measure of "goodness"

Cost function / loss function

Measure of how far the current estimate deviate from training data

Training input  $(\bar{x}_j, y_j)$

Compare  $f(\bar{x}_j)$  &  $y_j$  — Accumulate over all  $\bar{x}_j$

Simple difference  $f(\bar{x}_j) - y_j$  may cancel

Should at least take  $|f(\bar{x}_j) - y_j|$

Theoretically, squared error is best

$$\frac{1}{2} \sum_j (f(\bar{x}_j) - y_j)^2$$



Constant factor does not affect the nature of the cost

Repeatedly adjust  $\theta_i$  so that cost reduces

In Statistics - can explicitly calculate best  $\theta_i$  using mean, variance etc

How to iterate  $\theta_i$ ?

Initial estimate  $E_0 = \langle \theta_0^0, \theta_1^0, \dots, \theta_k^0 \rangle$

$\downarrow$   
 $C_0 \xrightarrow{\text{adjust}} E_1 \rightarrow C_1 \dots$

How to go from  $C_j$  to  $E_{j+1}$ ?

Gradient descent (greedy - reduce cost as much as possible)

For each parameter  $\theta_j$ , compute  $\frac{\partial C}{\partial \theta_j}$

Adjust  $\theta_j$  to  $\theta_j - \alpha \frac{\partial C}{\partial \theta_j}$

$\uparrow$   
step factor

Cost is  $\frac{1}{2} \sum_{k=1}^N (f(x_k) - y_k)^2$

$$f(\bar{x}_n) = \theta_0 + \theta_1 x_1^k + \theta_2 x_2^k + \dots + \theta_m x_m^k$$

$$\frac{\partial f(\bar{x}_n)}{\partial \theta_j} = \bar{x}_j^k$$

$$\frac{\partial}{\partial \theta_j} \frac{1}{2} (f(\bar{x}_n) - y)^2$$

$$\cancel{\frac{1}{2}} \cdot \cancel{2} (f(\bar{x}_n) - y)_k \cdot \left( \frac{\partial f(\bar{x}_n)}{\partial \theta_j} \right) x_j$$



$$\theta_j^{k+1} = \theta_j^k + \alpha \sum_{i=1}^l (y_i - f(\bar{x}_i)) \cdot \bar{x}_i$$

Finding optimal  $\theta_i$  — regression

Start with some initial estimate  $\theta_0$

Calculate  $f(\bar{x}_j)$  for all  $j$  & compute  
gradient for each coefficient

Update each  $\theta_j$  as above

Batch gradient descent

Update  $\theta_j$ 's after a full batch  
of predictions

Adjustment of  $\theta$  is proportional to  
error  $f(x_j) - y_j$

Another option

Update  $\theta_j$  after each input

# Stochastic Gradient Descent

- └ More unstable
- └ Overall converges faster

Overfitting

Regularization — keep coefficients simple

Add to cost, an extra penalty proportional to size of coefficients

Decision trees suffer from variance

Small input perturbation  $\rightarrow$  big change  
in output

Regression suffers from bias

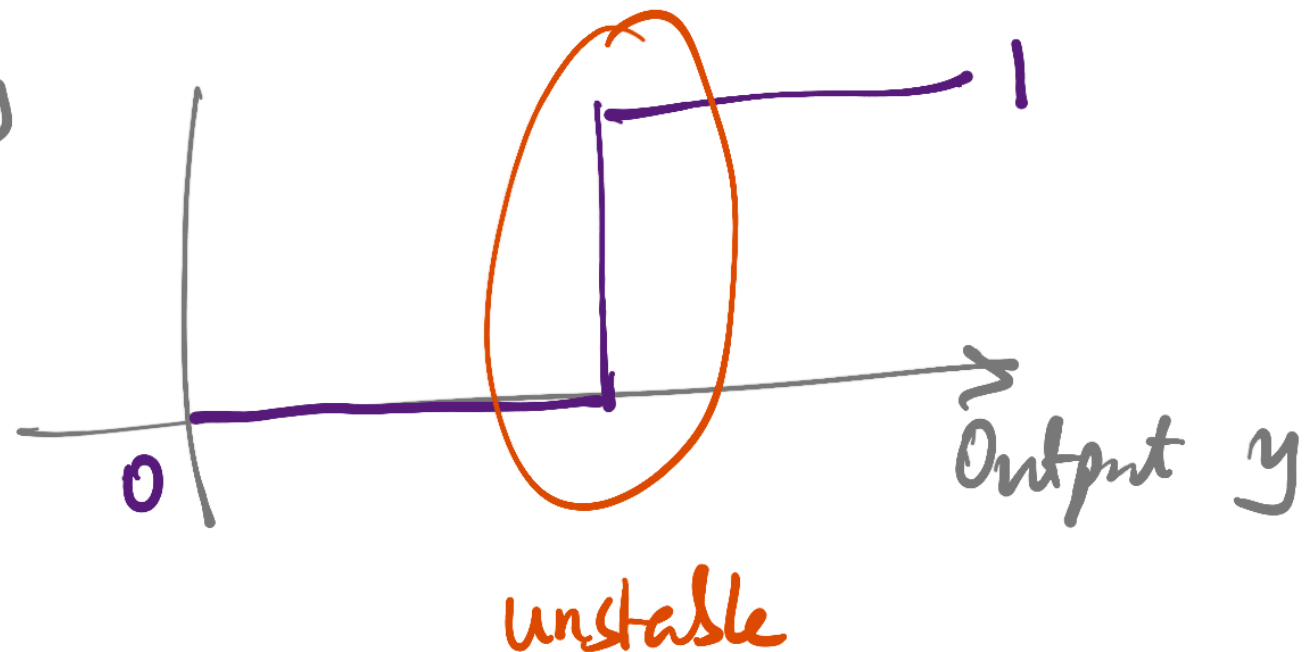
Intentionally limiting the nature of  $f(x)$

Model chosen is not expressive enough

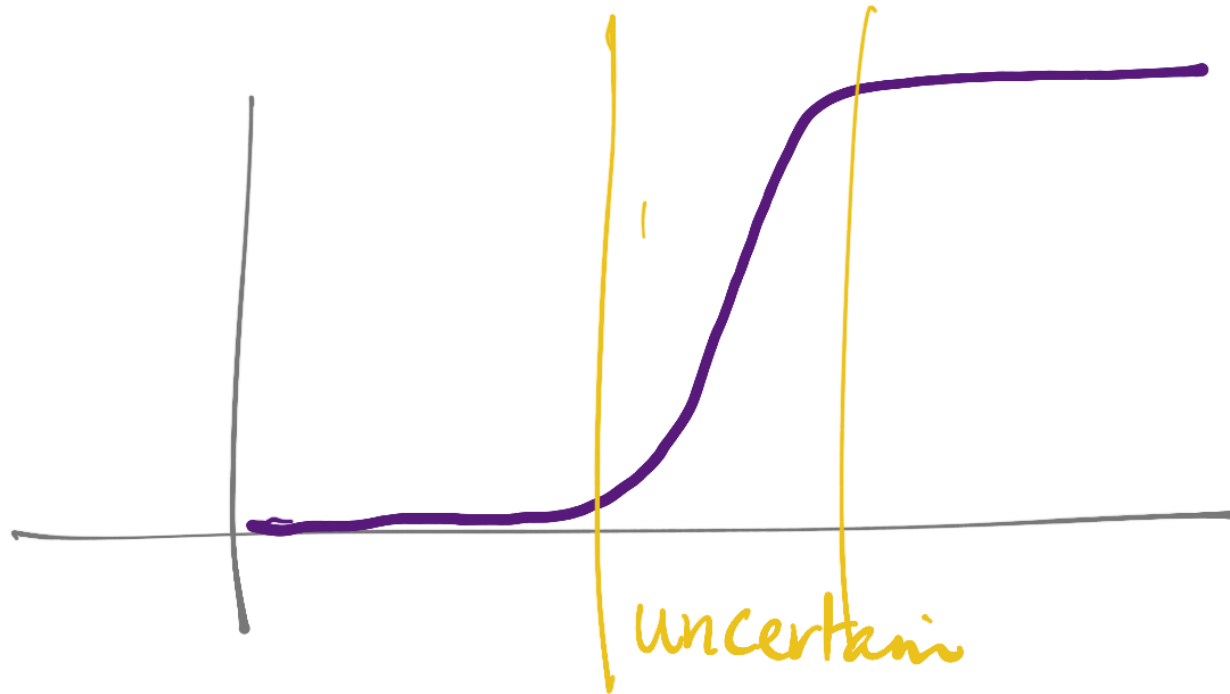
# Regression for classification

Classification  $\rightarrow$  predict 0 or 1

Typically



Instead, a smooth approximation of step



Sigmoid →

$$\sigma(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$[\theta_0 \ \theta_1 \ \dots \ \theta_m] \quad \bar{\theta}^T$$

$$\begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \quad \bar{x}$$

Still use square-error for cost

$$\frac{\partial C}{\partial \theta_j} = (f(x_j) - y_j) \cdot \frac{\partial f}{\partial \theta_j}$$

can be  
calculated

$$\sigma(z) = \frac{1}{1 + e^{-\theta z}}$$

$$\sigma'(z) = \sigma(z) (1 - \sigma(z))$$

Logistic regression



Sigmoid is not only option for smooth step

Square is not only option for cost

---

Regression using decision trees

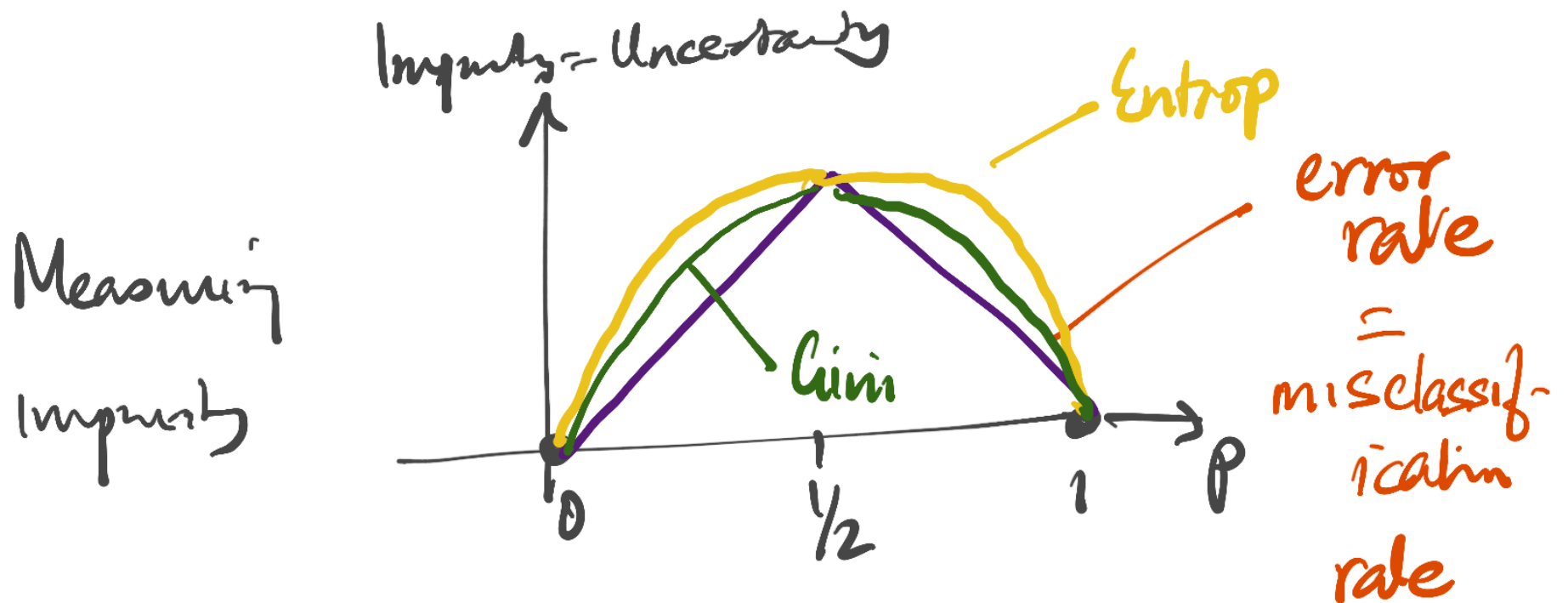
R, Python — CART

Classification & Regression Tree

# CART

- Classification - same as a decision tree

$1 - \sum p_i^2$  - Gini index instead of entropy

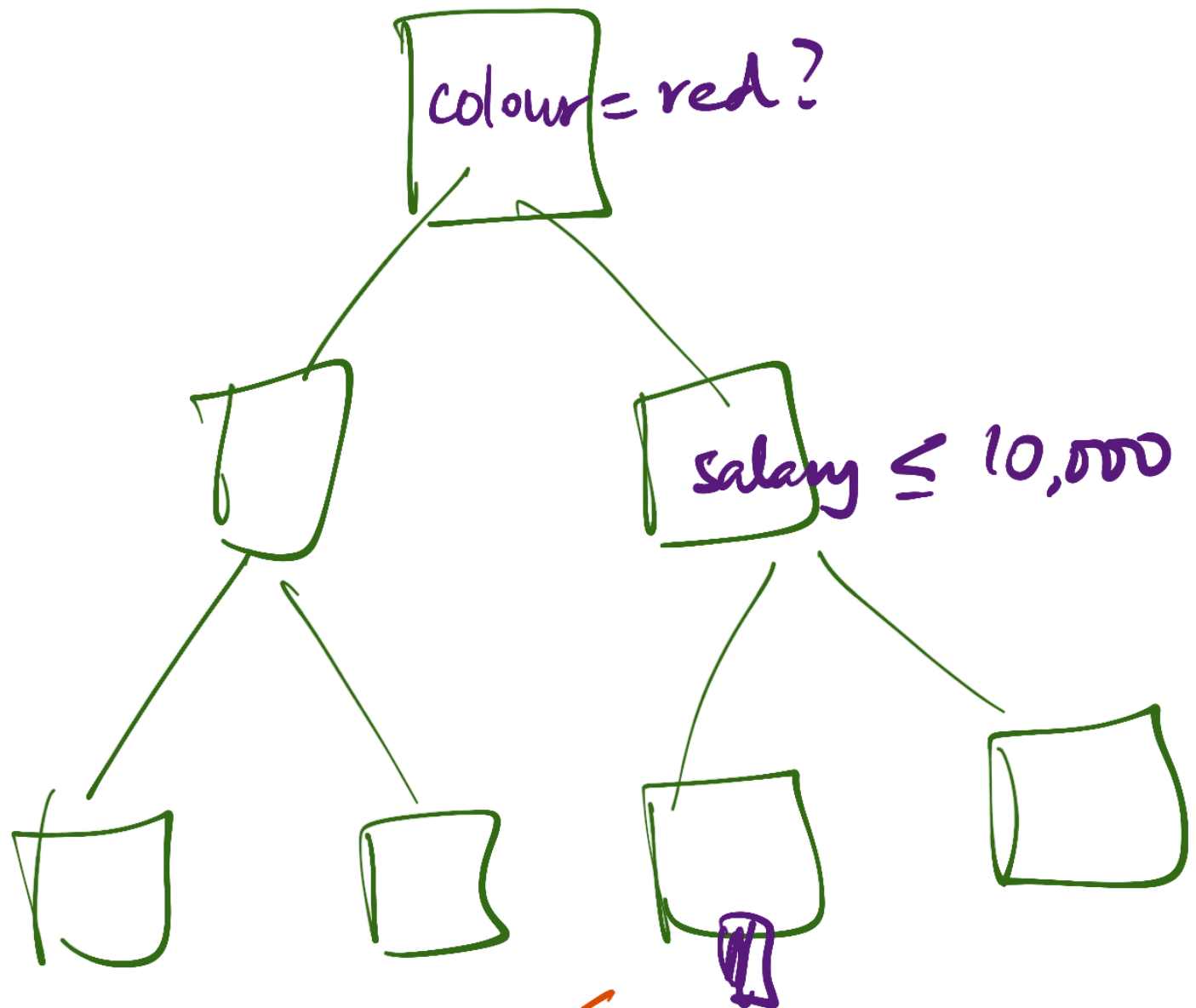


- Restrict to binary branching

## Regression trees

- ① - Which attribute/question to ask next
- ② - How to evaluate cost/purity
- ③ - What value to assign at a leaf

3



Take mean of  
all  $y$ 's in this node

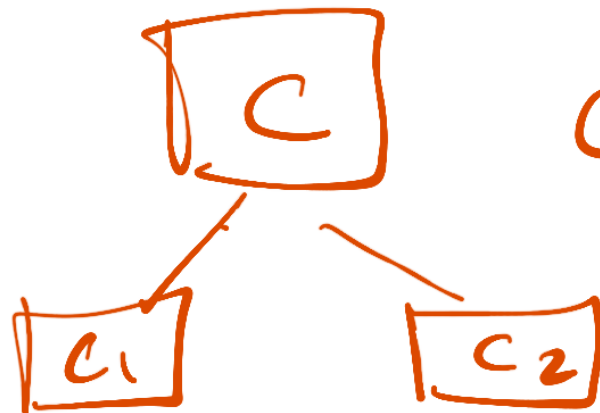
Predict  $y$ ?

② & ① Use the usual cost function

$f(x)$  is the average for that node

$$\sum (y_{\text{mean}} - y_i)^2$$

Choose question to reduce cost



$C = (C_1 + C_2)$  to be maximized