

## Decision Trees

### Predicting loan defaults with decision trees

What makes a loan risky?

"I want to buy"  
a new house

Loan Application

Credit  
\*\*\*

Income  
\*\*\*

Term.  
\*\*\*

Personal Info  
\*\* \*

Credit history

Did I pay previous  
loans on time?

Example: excellent  
good, or fair

What's my income?

Example: \$80k per year

Income

How soon do I need to  
pay the loan?

Example: 3 years,  
5 years

Term

Age, reason for the loan  
marital status, ...

Example: home loan  
for a married couple

Personal info

Intelligent application

Loan Applications

#1 → Intelligent loan

#2 → application

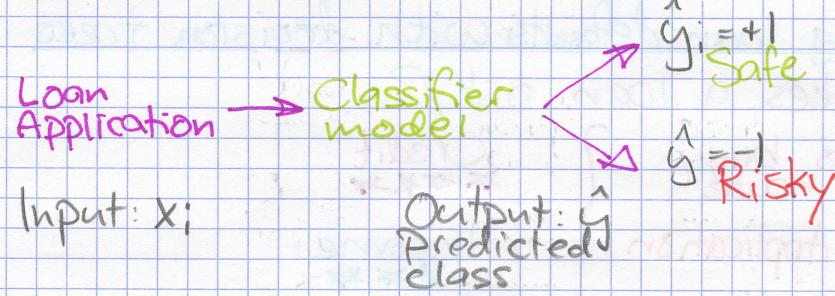
#3 → review system

Safe

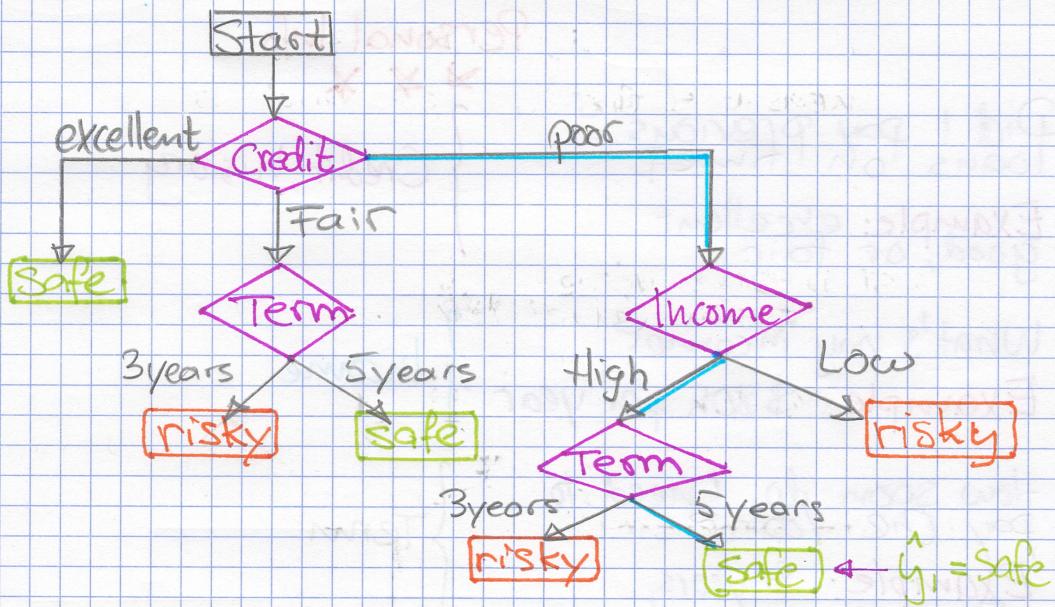
Risky

Risky

## Classifier review



This module ... decision trees



Scoring a loan application

$$x_i = (\text{Credit} = \text{poor}, \text{Income} = \text{high}, \text{Term} = 5 \text{ years})$$

Decision tree learning task

Learn decision tree from data?

Credit	Term	Income	y	Training data N observations ( $x_i, y_i$ )
excellent	2 yrs	high	safe	
fair	5 yrs	low	risky	
poor	3 yrs	high	safe	
poor	5 yrs	high	risky	
excellent	3 yrs	low	risky	
fair	5 yrs	low	safe	
poor	3 yrs	high	risky	
poor	5 yrs	low	safe	
fair	3 yrs	high	safe	

$h_1(x)$        $h_2(x)$        $h_3(x)$       loan status

Optimize quality metric on training data  $T(x)$

Quality metric: Classification error

- Error measures fraction of mistakes

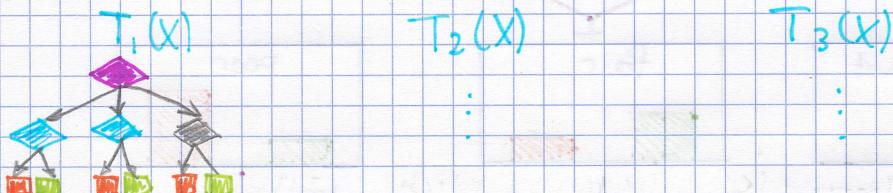
$$\text{Error} = \frac{\# \text{Incorrect predictions}}{\# \text{Examples}}$$

- Best possible value: 0%

- Worst possible value: 100%

How do we find the best tree?

Exponentially large number of possible trees makes decision tree learning hard!  
(NP-hard problem)

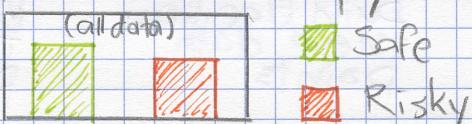


Simple (greedy) algorithm finds "good" trees

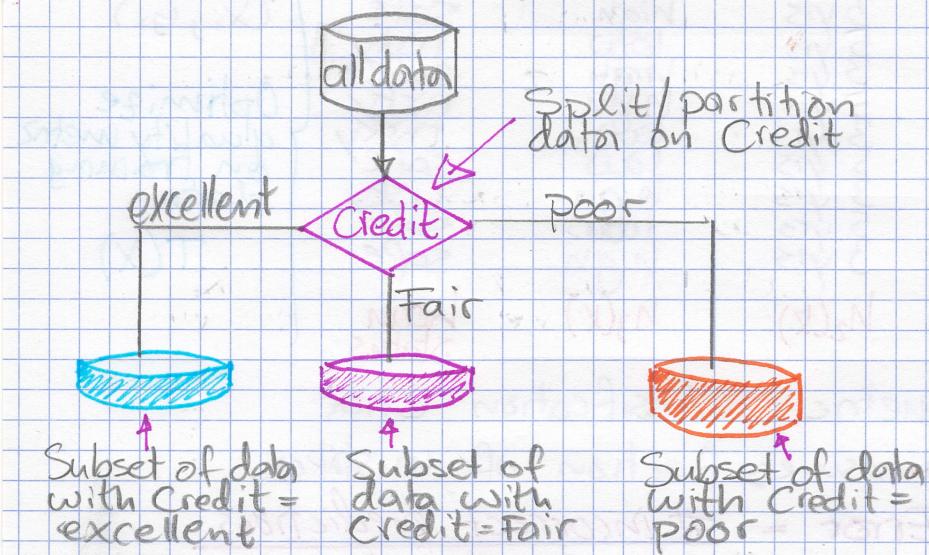
Greedy decision tree algorithm:

Algorithm outline

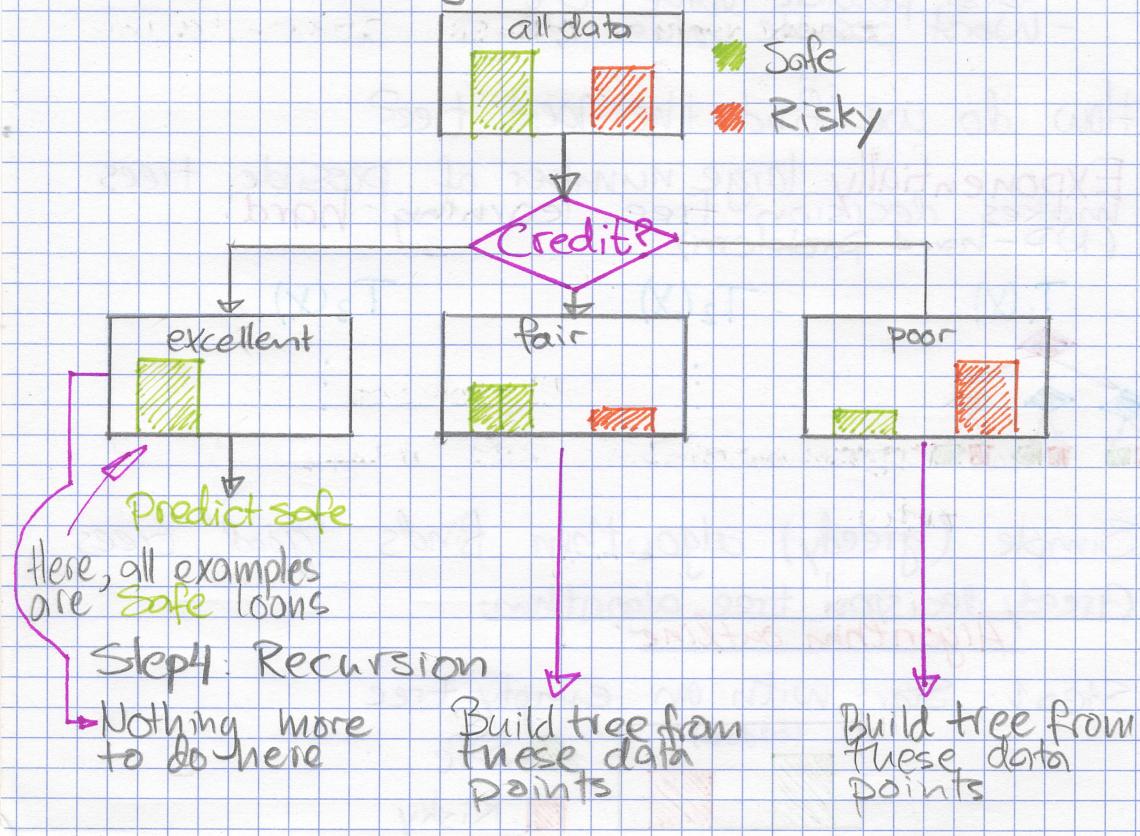
Step 1: Start with an empty tree



## Step 2: Split on a feature!



## Step 3: Making predictions



## Greedy decision tree learning

- Step 1: Start with an empty tree
- Step 2: Select a feature to split data
- For each split of the tree:
  - Step 3: If nothing more to make prediction
  - Step 4: Otherwise, go to step 2 & continue (recurse) on this split

Recursion

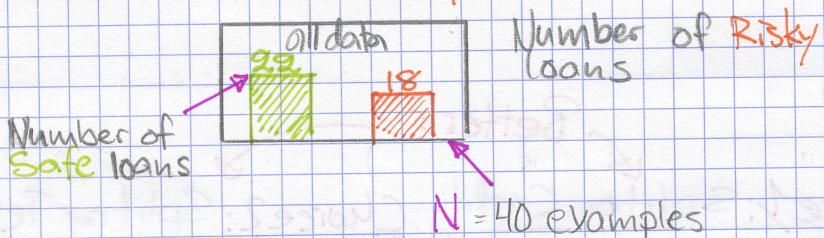
Problem 1: Feature split selection

Problem 2: Stopping condition

Learning a decision stump  
Feature split learning = decision stump learning

Start with all the data

Loan status Safe Risky

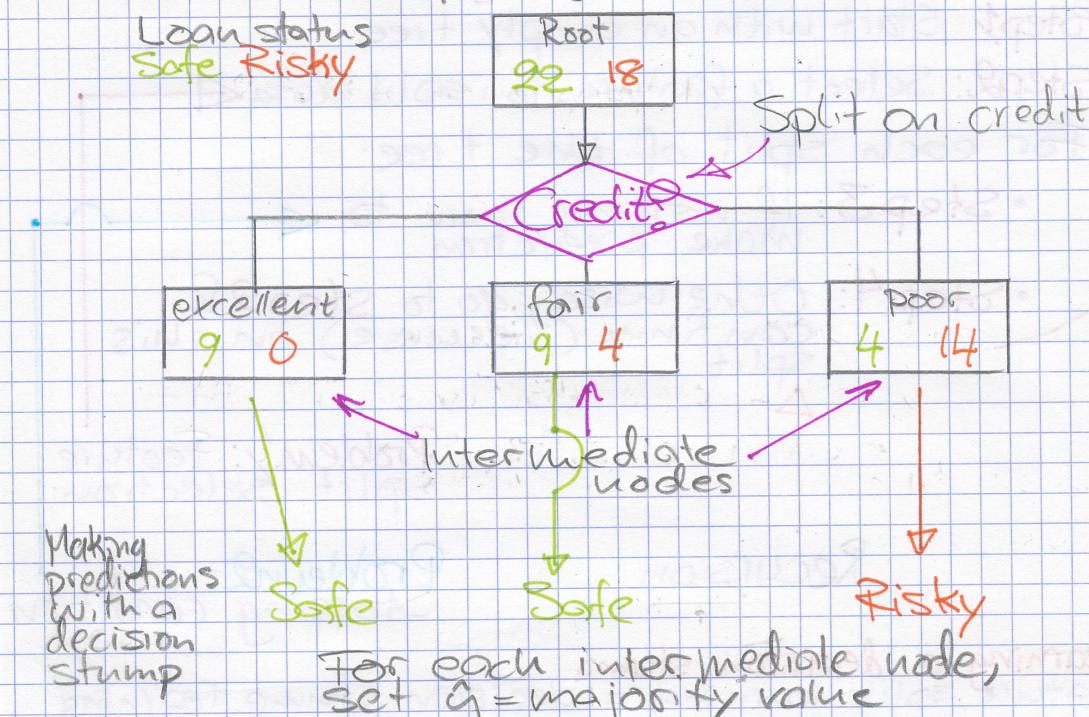


Compact visual notation: Root node

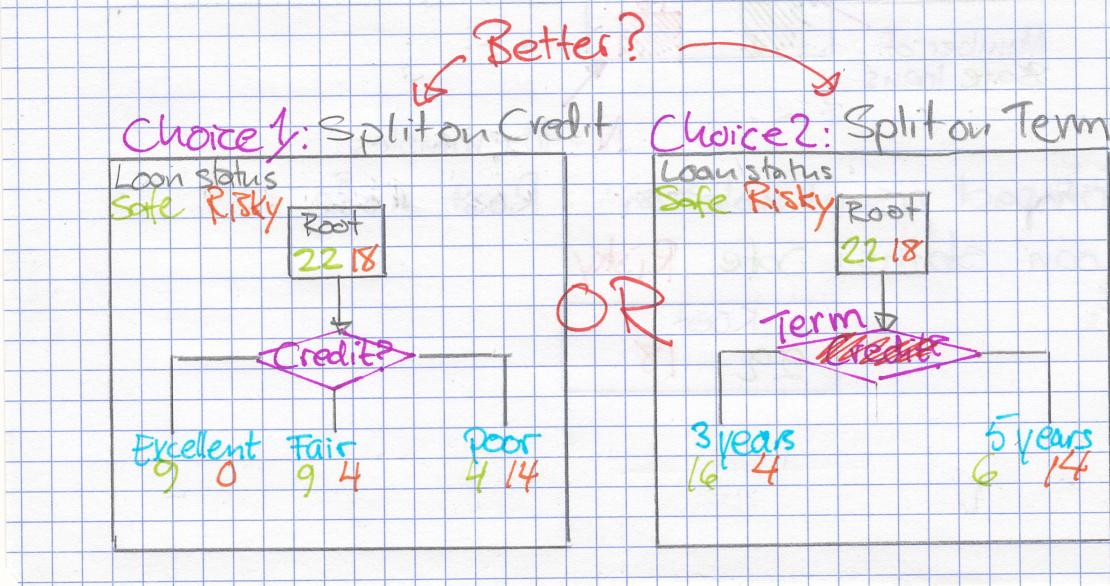
Loan status Safe Risky



## Decision Stump: Single level tree



Selecting best feature to split on  
How do we select the best feature



How do we measure effectiveness of a split?

Loan status  
safe Risky Root  
22 18

Credit

excellent fair  
9 0

Idea: Calculate classification error of this decision stump

poor  
9 4

4 14

$$\text{Error} = \frac{\# \text{mistakes}}{\# \text{data points}}$$

Calculating classification Error

- Step 1:  $\hat{y}$  = class of majority of data in node
- Step 2: Calculate classification error of predicting  $\hat{y}$  for this data

Loan status  
safe Risky

Root  
22 18

22 correct

Safe

18 mistakes

$$\text{Error} = \frac{18}{22+18} = 0.45$$

$\hat{y}$  = majority class

Tree	Classification error
(root)	0.45

Choice 1: Split on credit history

Loan stats  
Safe Risky

Root  
22 18

Credit?

excellent

9 0

Safe

↑

0 mistakes

fair

9 4

safe

↑

4 mistakes

poor

4 14

Risky

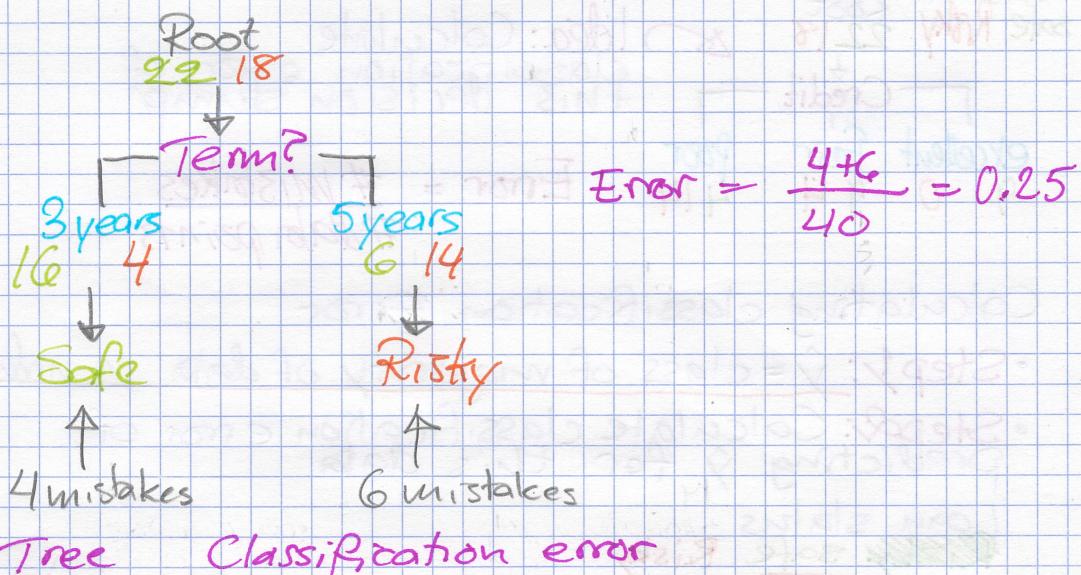
↑

4 mistakes

Does a split on Credit reduce classification error below 0.45?

$$\text{Error} = \frac{4+4}{40} = 0.20$$

## Choice 2: Split on Term



(root)	0.45	
<u>Split on Credit</u>	0.2	First split
<u>Split on Term</u>	0.25	

Choice 1 vs Choice 2

Feature split selection algorithm

- Given a subset of data  $M$  (a node in a tree)
- For each feature  $h_i(x)$ : credit, term, income
- 1: Split data of  $M$  according to feature  $h_i(x)$
- 2: Compute classification error split
- Choose feature  $h^*(x)$  with lowest classification error

## When to stop recursing

We've learned a decision stump, what next?

Loan status  
Safe Risky

Root  
22 18

Credit

excellent  
9 0

fair  
9 4

poor  
4 14

Safe

Leaf node

All data points are  
**Safe** → nothing  
else to do with  
this subset of  
data

Build decision  
stump with  
subset of data  
where Credit=fair

Build decision  
stump with  
subset of data  
where Credit=poor

## Simple greedy decision tree learning

Pick best feature to split on

↑ Learn decision stump with  
this split

For each leaf of decision  
stump, recurse

When do we stop?

Stopping condition 1: All data agrees on y

All data in these  
nodes have same y value  
→ nothing to do

3 years  
0 2

excellent  
9 0

Stopping condition 2: Already split on all features

### Greedy decision tree learning

• Step 1: Start with an empty tree

• Step 2: Select a feature to split data

• For each split of the tree

• Step 3: If nothing more to make predictions

• Step 4: Otherwise, go to Step 2 & continue (recurse) on this split

Pick feature split leading to lowest classification error

Stopping conditions 1 & 2

Recursion

### Decision tree prediction algorithm

**predict**(tree-node, input)

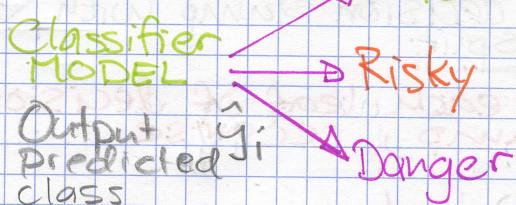
• If current tree-node is a leaf:  
o return majority class of data points in leaf

• else:  
o next-node = child of tree-node whose feature value agrees with input  
o return predict(next-node, input)

### Multiclass classification & predicting probabilities

#### Multiclass prediction

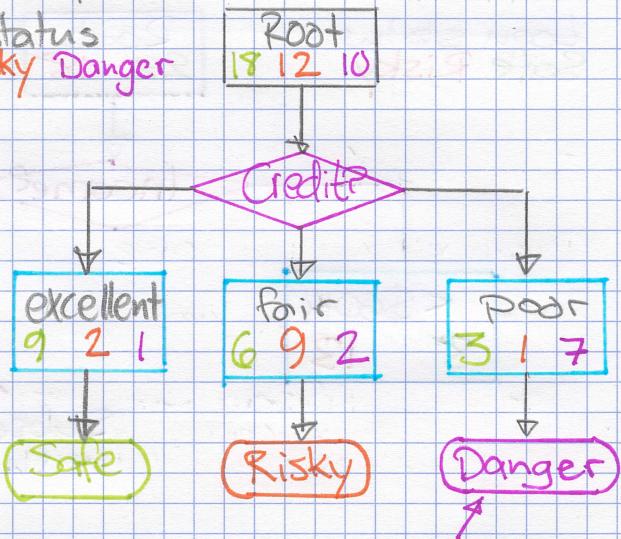
Loan Application →  
Input:  $x_i$



## Multiclass decision stump

$N = 40$   
1 feature  
3 classes

Loan status  
Safe Risky Danger



$$P(y = \text{danger} | x)$$

$$= \frac{7}{3+1+7} = 0.64$$

Decision tree learning:  
Real valued features

How do we use real values inputs?

Income	Credit	Tenn	y
\$105k	excellent	3 yrs	Safe
112	good	5 yrs	Risky
73	fair	3 yrs	Safe
69	excellent	5 yrs	Safe
217	Excellent	3 yrs	Risky
120	good	5 yrs	Safe
64	fair	3 yrs	Risky
340	excellent	5 yrs	Safe
60	good	3 yrs	Risky

Split on each numeric value?

Danger: May  
only contain  
one data point  
per node

Root: 22.18  
Loan status  
Safe Risky

Income

\$30K 31.4K 39.5K ----- \$61.1K \$91.3K

Can't trust prediction (overfitting)

## Alternative: Threshold split

Loan status  
Safe Risky

Root  
22 18

Split on the feature income

Income

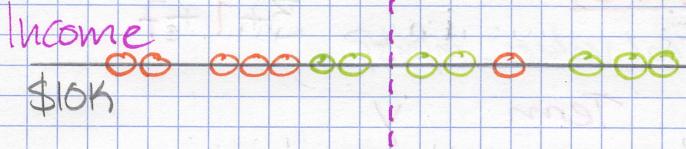
< \$60K  
8 13

$\geq \$60K$   
14 5

Subset of data with  $\text{income} \geq \$60K$

Many data points  $\rightarrow$  lower chance of overfitting

Threshold splits in 1-D  
Threshold split is the line  $\text{income} = \$60K$



Safe  
Risky

Each split partitions the 2-D space

