

Handling Missing Data

So far: data always completely observed

"Known x and y values for all data points"

Missing data

Credit	Term	Income	y
excellent	3 yrs	high	safe
fair	?	low	risky
fair	3 yrs	high	safe
poor	5 yrs	high	risky
excellent	3 yrs	low	risky
fair	5 yrs	high	safe
poor	?	high	risky
poor	5 yrs	low	safe
fair	?	high	safe

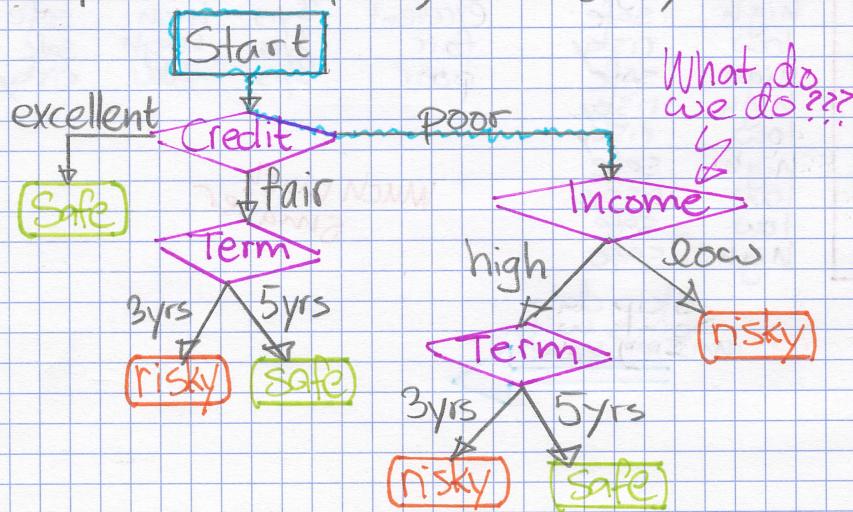
← Loan application may be 3 or 5 years

Missing values impact training and predictions

- 1 Training data: Contains "unknown" values
- 2 Predictions: Input at prediction time contains "unknown" values

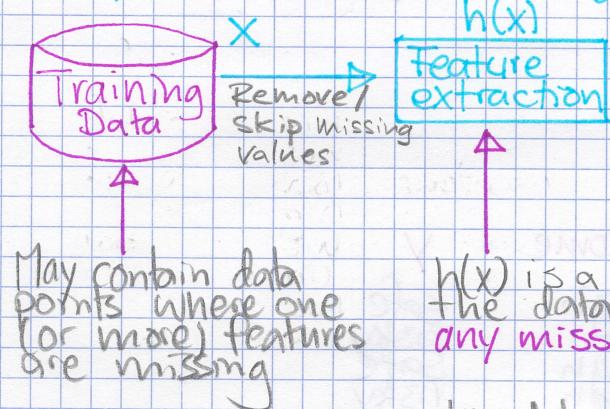
Missing values during prediction

$x_i = (\text{credit}=\text{poor}, \text{income}=?, \text{term}=5 \text{ years})$



Strategy 1: Purification by skipping missing data

Ideally: Purification by skipping / removing



May contain data points where one (or more) features are missing

$N=9$, 3 features

Credit Term Income y
excellent 3yrs high safe
:
:

Skip data points with missing values

$N=6$, 3 features

Credit Term Income y
:
:

The challenge with Ideal

Warning: More than 50% of the loan Terms are unknown

N=9, 3 features			
Credit	Term	Income	y
excellent	3yrs	high	safe
fair	?	low	risky
fair	3yrs	high	safe
poor	?	high	risky
excellent	?	low	risky
fair	?	high	safe
poor	3yrs	low	risky
poor	?	low	safe
fair	?	high	safe

N=3, 3 features			
Credit	Term	Income	y
excellent	3yrs	high	safe
fair	3yrs	high	safe
poor	3yrs	low	risky

Much much smaller

Skip data points with missing values

Idea 2: Skip features with missing values

Skip features with many missing values

$h(x)$	N=9, 2 features	Credit	Income	y
		excellent	high	safe
		fair	low	risky
		fair	high	safe
		poor	high	risky
		excellent	low	risky
		fair	high	safe
		poor	high	risky
		poor	low	safe
		fair	high	safe

Missing value skipping: Ideas 1 & 2

Idea 1: Skip data points where any feature contains a missing value

- Make sure only a few data points are skipped

Idea 2: Skip an entire feature if it's missing for many data points

- Make sure only a few features are skipped.

Missing value skipping: Pros and Cons

Pros

- Easy to understand and implement
- Can be applied to any model (decision trees, logistic regression, linear regression, ...)

Cons

- Removing data points and features may remove important information from data
- Unclear when it's better to remove data points versus features
- Doesn't help if data is missing at prediction time

Strategy 2: Purification by imputing missing data

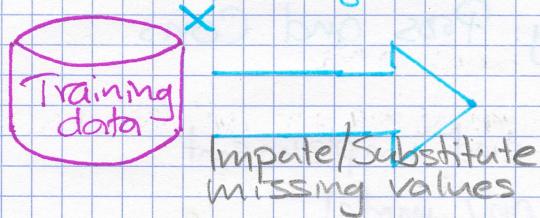
Main drawback of skipping strategy

X				h(x)			
N=9, 3 features				N=6, 3 features			
Credit	Term	Income	y	Credit	Term	Income	y
excellent	3 yrs	high	safe	excellent	3 yrs	high	safe
fair	?	low	risky	fair	3 yrs	high	safe
fair	3 yrs	high	safe	poor	5 yrs	high	risky
poor	5 yrs	high	risky	excellent	3 yrs	low	risky
excellent	3 yrs	low	risky	fair	5 yrs	high	safe
fair	5 yrs	high	safe	poor	3 yrs	low	risky
poor	3 yrs	low	risky				
poor	?	low	safe				
fair	?	high	safe				

Skip data points with missing values

Can we keep all the data?

Idea 2: Purification by imputing



Use other data points in X to "guess" the "?"

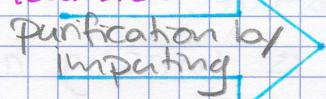
h(x)

Feature extraction

Same number of data points as the original data

Example: Replace ? with most common value

3 year loans: 4 ← Best guess
5 year loans: 2



Credit	Term	Income	y
excellent	3 yrs	high	safe
fair	3 yrs	low	risky
fair	3 yrs	high	safe
Poor	5 yrs	high	risky
excellent	3 yrs	low	risky
fair	5 yrs	high	safe
Poor	3 yrs	low	risky
Poor	3 yrs	low	safe
Fair	3 yrs	high	safe

Common (simple) rules for purification by imputation

Impose each feature with missing values:

1) Categorical features use mode: Most popular value (mode) of non-missing x_i

2) Numerical features use average or median: Average or median value of non-missing x_i

Many advanced methods exist, e.g. expectation-maximization (EM) algorithm

Missing Value Imputation: Pros and Cons

Pros

- Easy to understand and implement
- Can be applied to any model (decision trees, logistic regression, linear regression)
- Can be used at prediction time: use same imputation rules

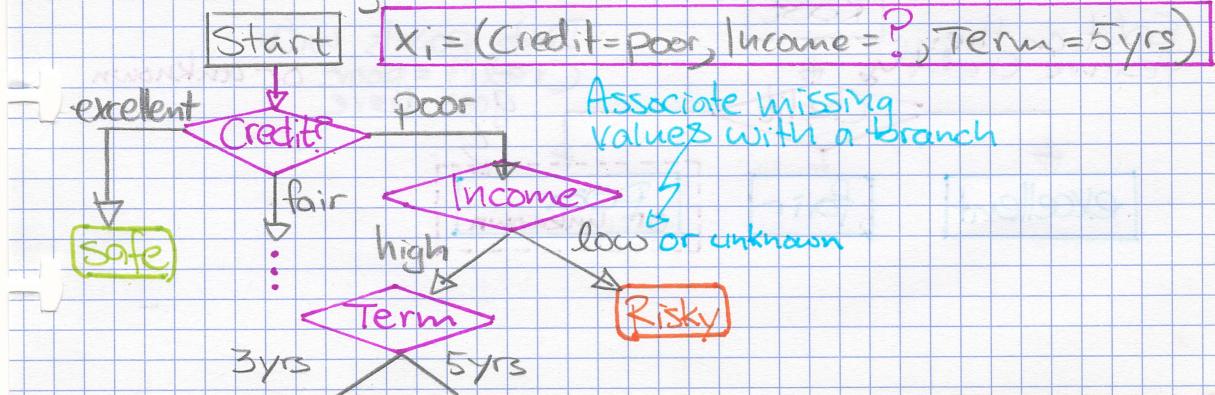
Cons

- May result in systematic errors

Example: Feature "age" missing in all banks in Washington by state law

Strategy 3: Adapt learning algorithm to be robust to missing values.

Add missing values to the tree definition



Add missing value choice to every decision node

Explicitly handling missing data by learning algorithm: Pros and Cons

Pros

- Addresses training and prediction time
- More accurate predictions

Cons

- Requires modification of learning algorithm
 - Very simple for decision trees

Feature split selection with missing data

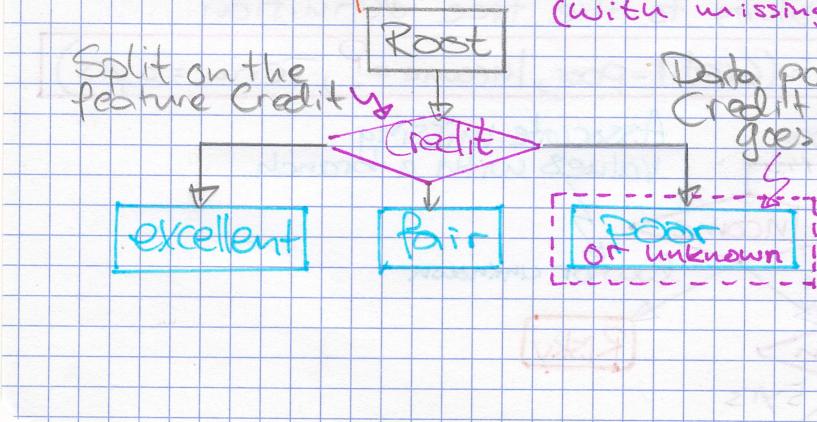
Greedy decision tree learning

- 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

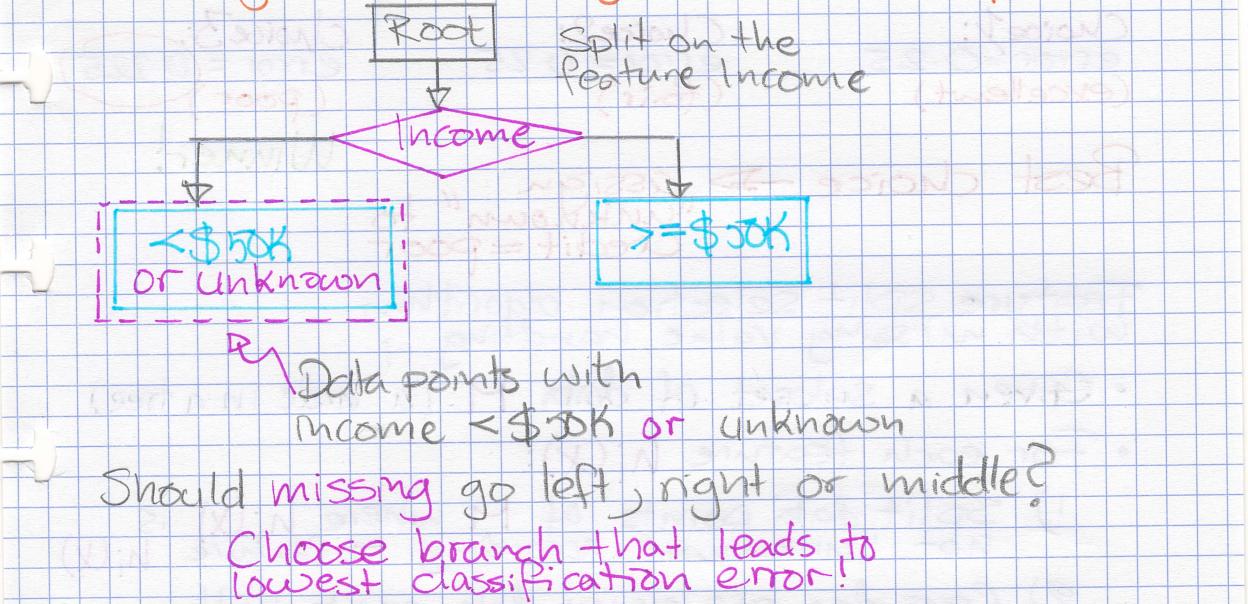
Pick feature split leading to lowest classification error

Must select feature & branch for missing values!

Feature split (without missing values) (with missing values)



Missing value handling in threshold splits



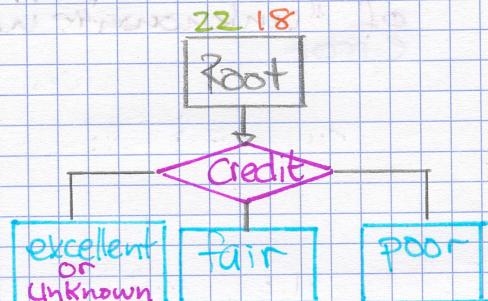
Choice 1: Missing values go with Credit = excellent

Choice 2: Missing values go with Credit = fair

Choice 3: Missing values go with Credit = poor

Computing classification error of decision stump with missing data

Credit	Term	Income	Y
Excellent	3 yrs	high	Safe
?	5 yrs	low	risky
Fair	3 yrs	high	Safe
?	5 yrs	high	risky
Poor	3 yrs	low	risky
?	5 yrs	low	safe
Poor	3 yrs	high	risky
Poor	5 yrs	low	safe
Fair	3 yrs	high	Safe
...



$$\begin{aligned}
 \text{Error} &= \frac{2+4+4}{40} \\
 &= 0.25
 \end{aligned}$$

Observed values "Unknown"

9	0	8	4	4	12
1	2				
Total		10	②	8 ④	④ 12

Use classification error to decide

Choice 1:
error = 0.25
(excellent)

Choice 2:
error = 0.25
(fair)

Choice 3:
error = 0.225
(poor)

Winner!

Best choice \rightarrow Assign "unknown" to credit = poor

Feature split selection algorithm with missing value handling

- Given a subset of data M (a node in a tree)
- For each feature $h_i(x)$:
 - Split data points of M where $h_i(x)$ is not "unknown" according to feature $h_i(x)$
 - Consider assigning data points with "unknown" value for $h_i(x)$ to each branch
 - Compute classification error split & branch assignment of "unknown" values
- Choose feature $h^*(x)$ & branch assignment of "unknown" with lowest classification error