

Decision Trees

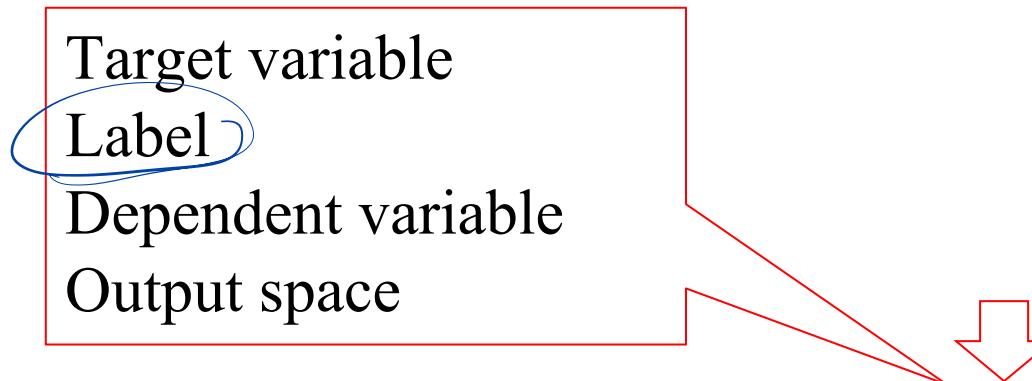
Nearest neighbor classification

MGTF 495

Class Outline

- Decision Trees
- Nearest Neighbor Classification
- Algorithmic Analysis
- Ensemble Methods
- Evaluation Metrics

Decision Trees

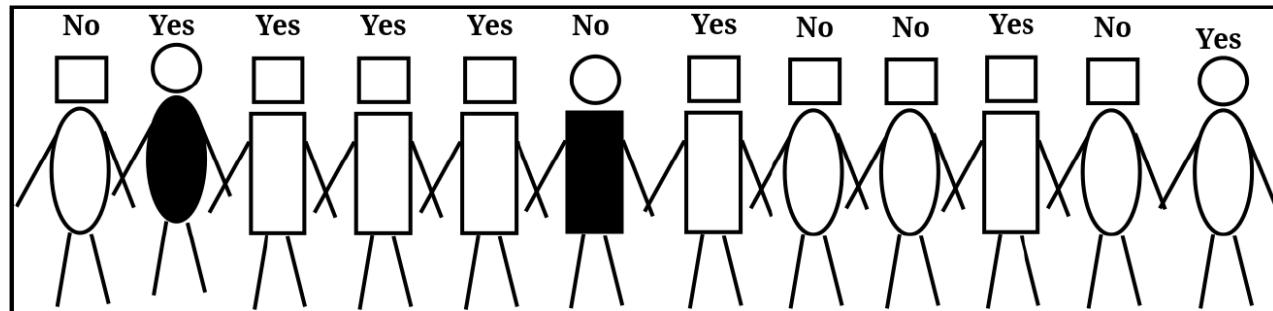


Person ID	Age	Gender	Income	Balance	Mortgage payment
123213	32	F	25000	32000	Y
17824	49	M	12000	-3000	N
232897	60	F	8000	1000	Y
288822	28	M	9000	3000	Y
....

Decision Trees

- How can we judge whether a variable contains important information about the target variable?
- How can we (automatically) obtain a selection of the more informative variables with respect to predicting the value of the target variable?
- Even better, can we obtain the ranking of the variables?

Decision Trees



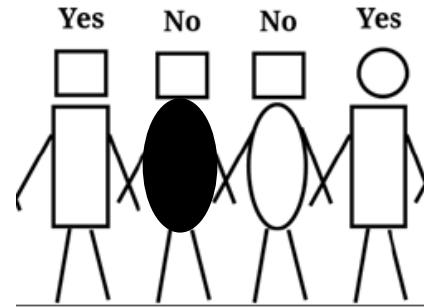
- Attributes:
 - head-shape: square, circular
 - body-shape: rectangular, oval
 - body-color: black, white
- Target variable: Yes, No

Decision Trees

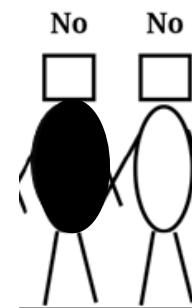
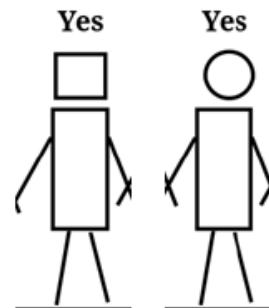
- Which attribute is the most informative? Or the most useful for distinguishing between data instances?
- If we split our data according to this variable, we would like the resulting groups to be as *pure* as possible.
- By pure we mean *homogeneous with respect to the target variable*.
- If every member of a group has the same value for the target, then the group is totally pure.

Example

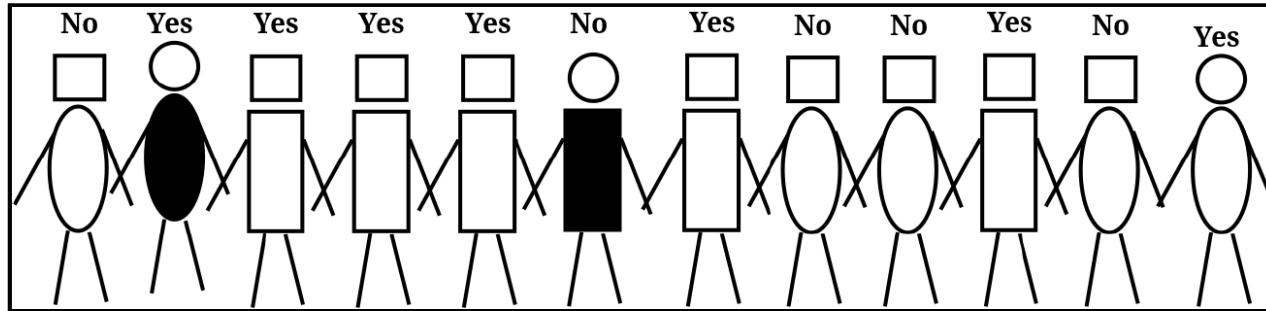
- If this is our entire dataset:



- Then, we can obtain two pure groups by splitting according to body shape:



Concerns



- Attributes rarely split a group perfectly.
- Even if one subgroup happens to be pure, the other may not.
- Is a very small, pure group, a good thing? *overfit*
- How should continuous and categorical attributes be handled?

Entropy and Information Gain

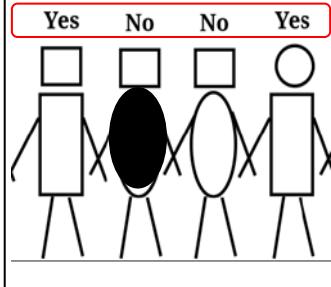
- Target variable has two (or more) categories: 1, 2 (...m)
- Probability P1 for category 1
- Probability P2 for category 2
- ...
- Entropy:

$$H(X) = -p_1 \log_2 p_1 - p_2 \log_2 p_2 - \dots - p_m \log_2 p_m$$

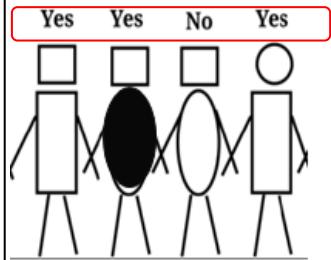
prob of seeing one value of this group

Entropy

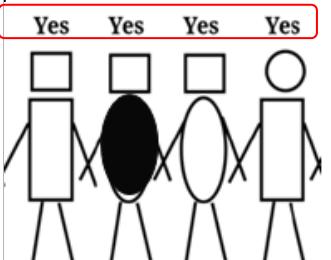
$$H(X) = -p_1 \log_2 p_1 - p_2 \log_2 p_2 - \dots - p_m \log_2 p_m$$



$$H(X) = -0.5 \log_2 0.5 - 0.5 \log_2 0.5 = 1$$



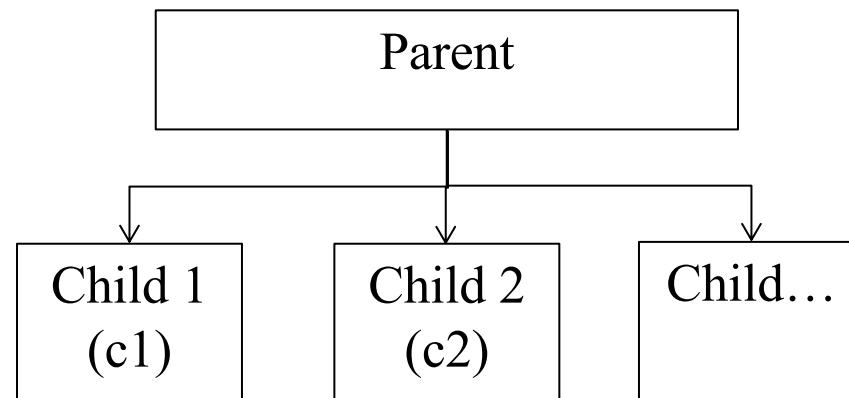
$$H(X) = -0.75 \log_2 0.75 - 0.25 \log_2 0.25 = 0.81$$



$$H(X) = -1 \log_2 1 = 0$$

Information Gain

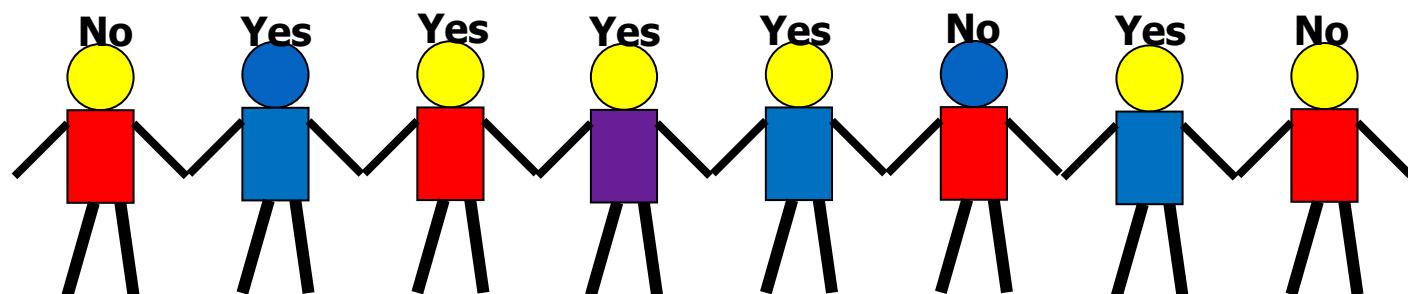
- Calculation of information gain (IG):
- $IG(\text{parent}, \text{children}) = \text{entropy}(\text{parent}) - [p(c_1) \times \text{entropy}(c_1) + p(c_2) \times \text{entropy}(c_2) + \dots]$



Note: Higher IG indicates a more informative split by the variable.

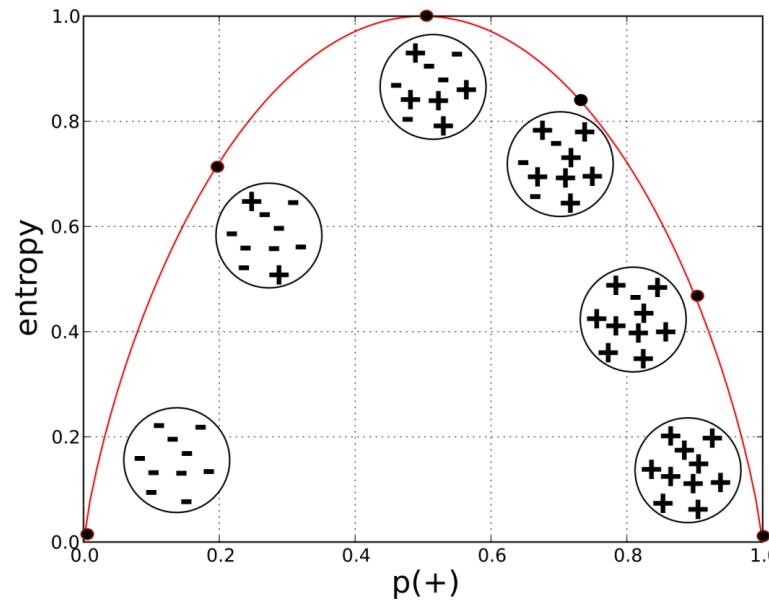
Selecting Informative Attributes

- Objective: Based on customer attributes, partition the customers into subgroups that are less impure – with respect to the class (i.e., such that in each group as many instances as possible belong to the same class)



Selecting Informative Attributes

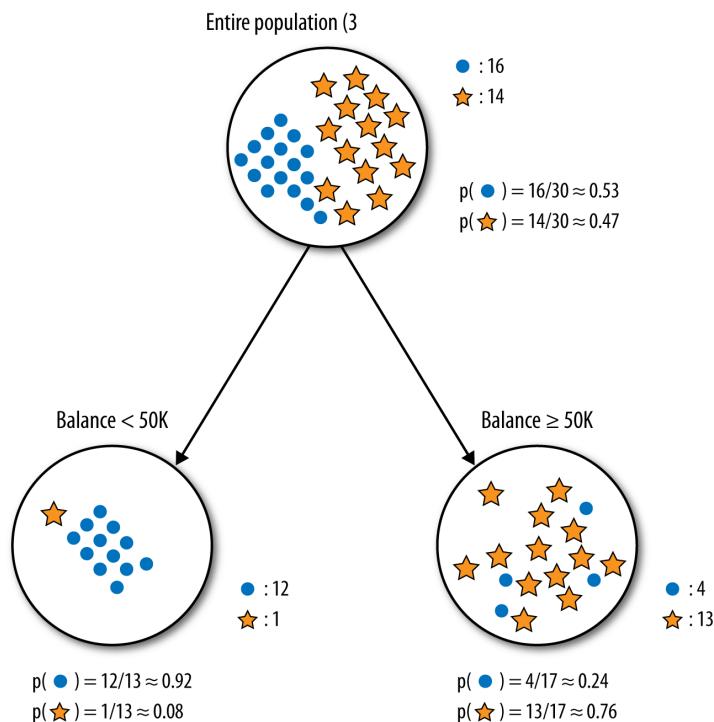
- The most common splitting criterion is called **information gain (IG)**
 - It is based on a **purity measure** called **entropy**
 - $entropy = -p_1 \log_2(p_1) - p_2 \log_2(p_2) - \dots$
 - Measures the general disorder of a set



Information Gain

- Information gain measures the change in entropy due to any amount of new information being added

$$IG(\text{parent}, \text{children}) = \text{entropy}(\text{parent}) - [p(c_1) \times \text{entropy}(c_1) + p(c_2) \times \text{entropy}(c_2) + \dots]$$



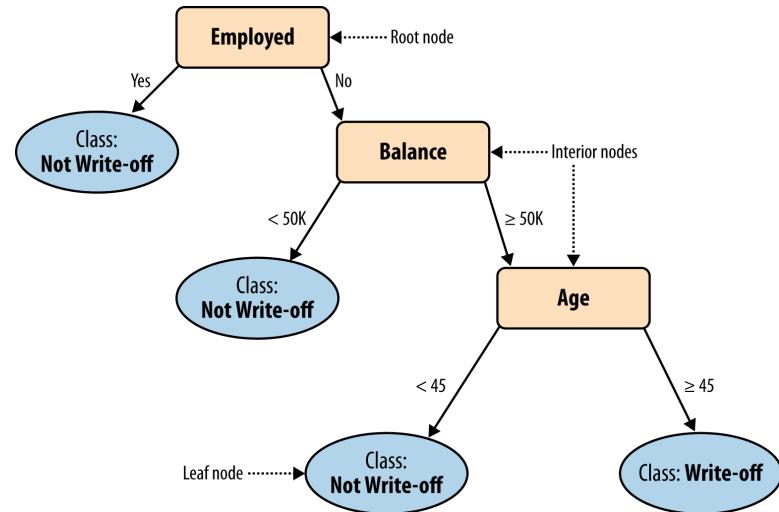
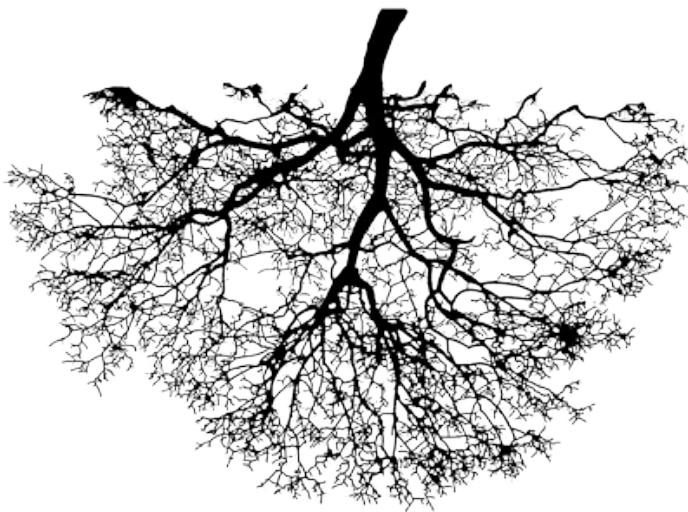
So far...

- We have measures of:
 - Purity of the data (entropy)
 - How informative is (a split by) a variable
- We can identify and rank informative variables
- Next – we will use this method to build our first supervised learning classifier – a decision tree

Multivariate Supervised Segmentation

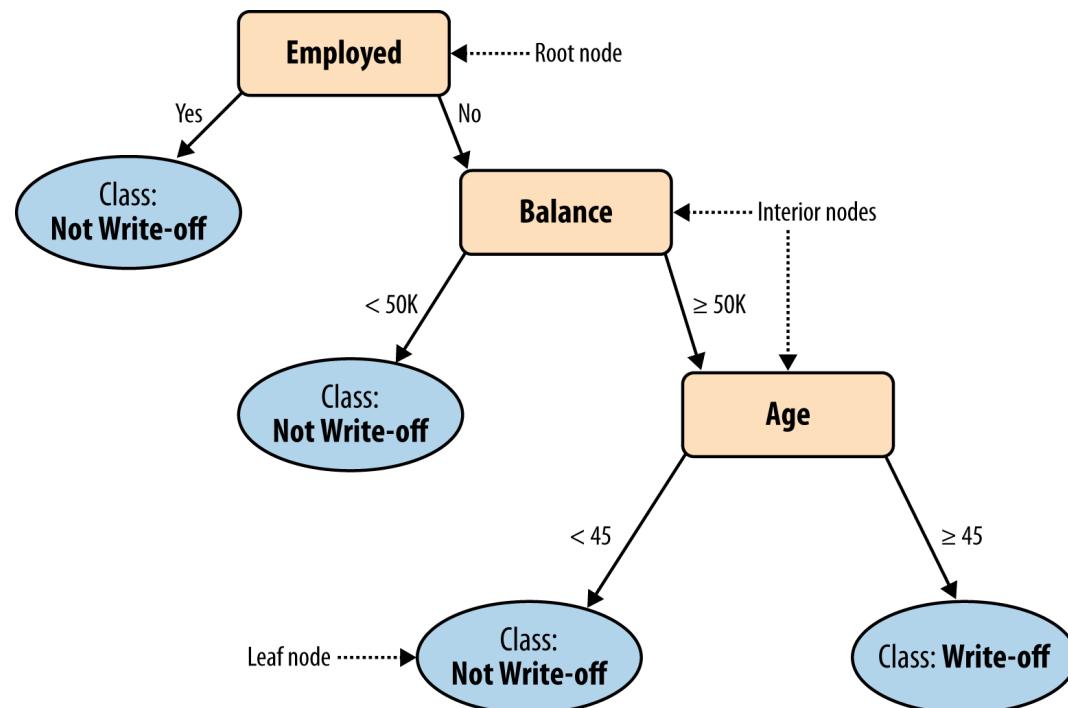
- If we select the *single* variable that gives the most information gain, we create a very *simple* segmentation
- If we select multiple attributes each giving some information gain, how do we put them together?

Tree-Structured Models



Tree-Structured Models

- Classify 'John Doe'
 - Balance=115K, Employed=No, and Age=40



Tree-Structured Models: “Rules”

- No two parents share descendants
- There are no cycles
- The branches always “point downwards”
- Every example always ends up at a leaf node with some specific class determination
 - Probability estimation trees, regression trees (*to be continued..*)

Tree Induction

- How do we create a classification tree from data?
 - **divide-and-conquer** approach
逐层递归
 - take each data subset and **recursively** apply attribute selection to find the best attribute to partition it
- When do we stop?
 - The nodes are pure,
 - there are no more variables, or
 - even earlier (over-fitting – *to be continued..*)

Why trees?

- Decision trees (DTs), or classification trees, are one of the most popular data mining tools
 - (along with linear and logistic regression)
- They are:
 - Easy to understand
 - Easy to implement
 - Easy to use
 - Computationally cheap
- Almost all data mining packages include DTs
- They have advantages for model comprehensibility, which is important for:
 - model evaluation
 - communication to non-DM-savvy stakeholders

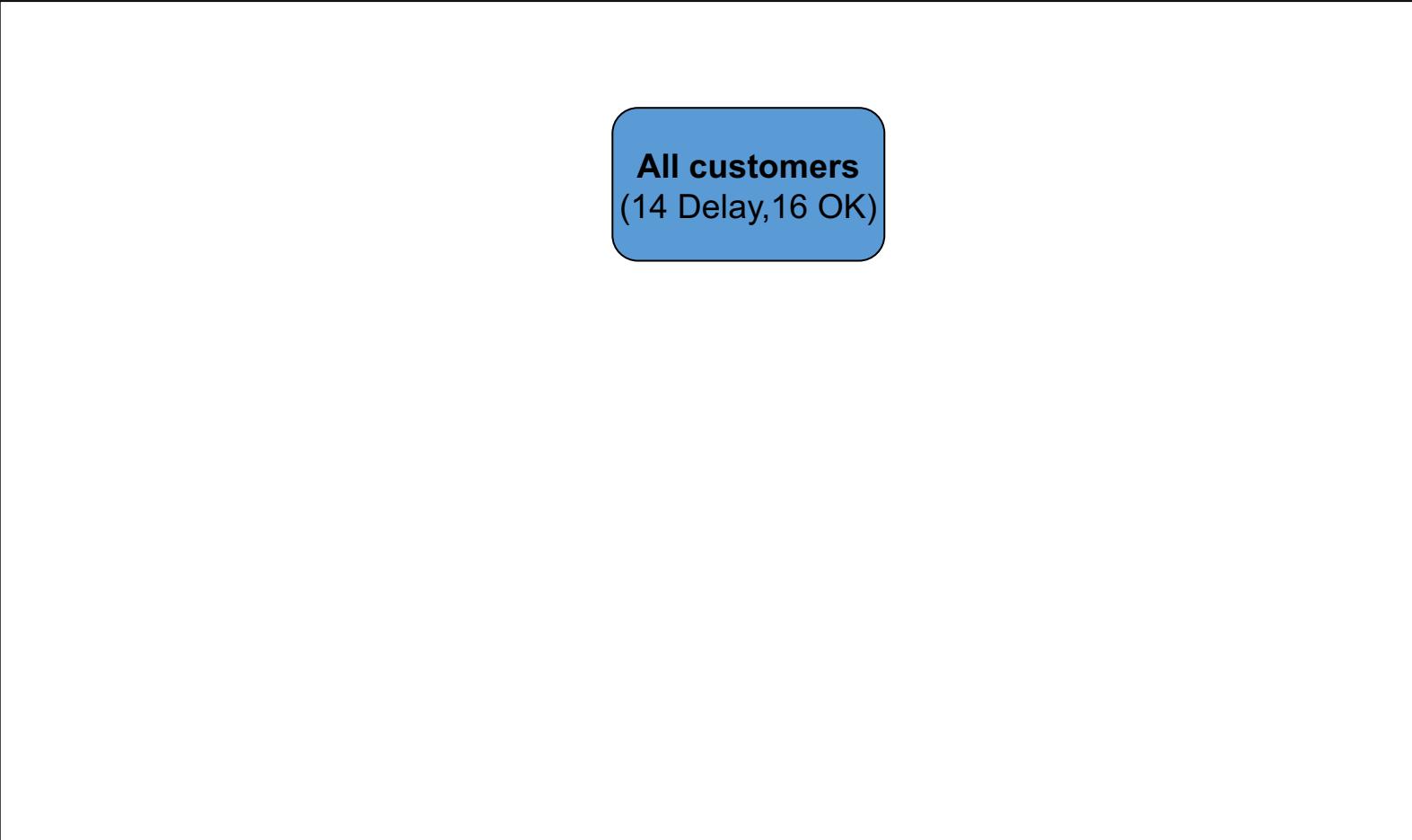
Dataset

Result

person id	age>50	gender	residence	Balance>= 50,000	mortgage payment delay
123213	N	F	own	N	delayed
17824	Y	M	own	Y	OK
232897	N	F	rent	N	delayed
288822	Y	M	other	N	delayed
....

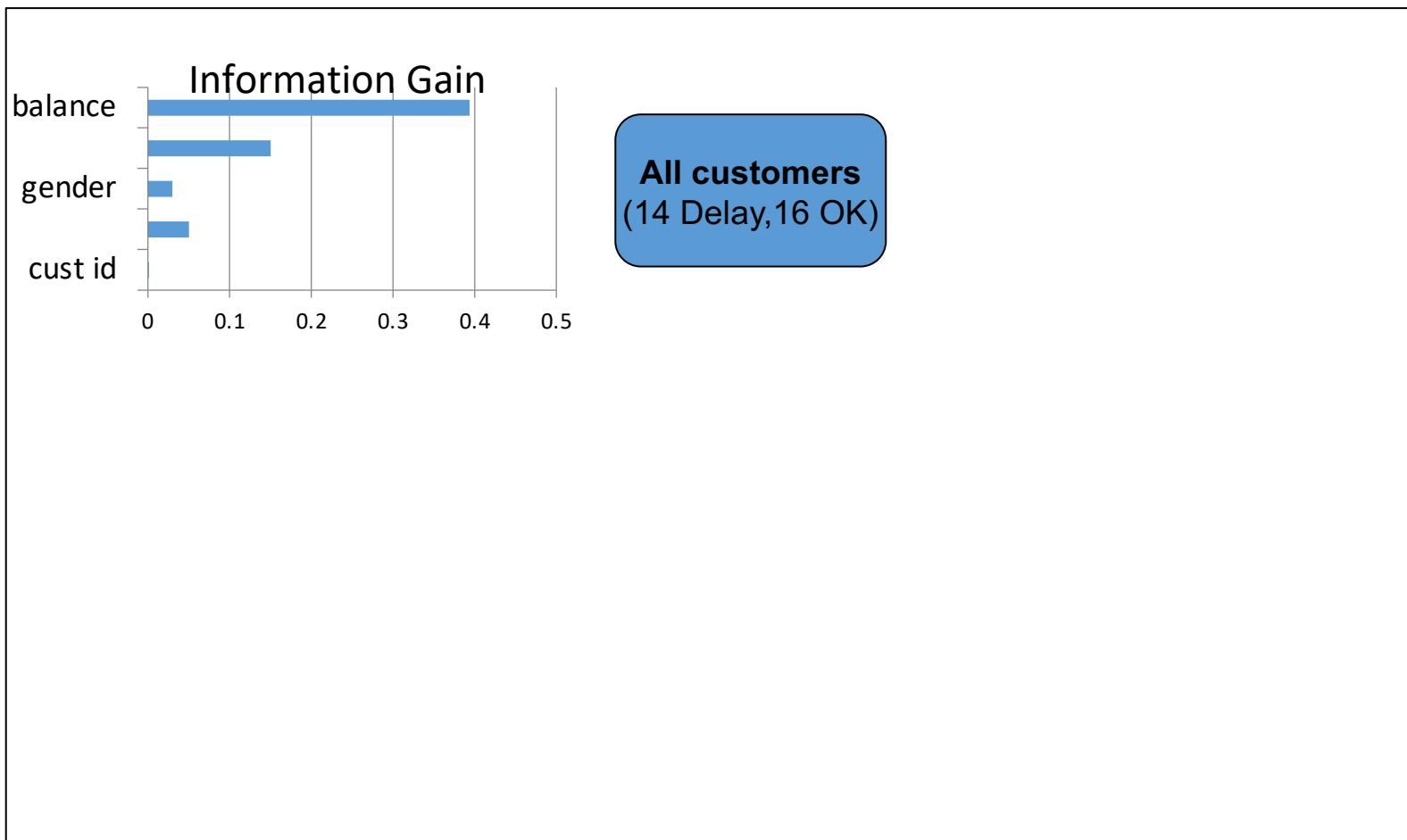
Based on this dataset we will build a tree-based classifier.

Tree Structure

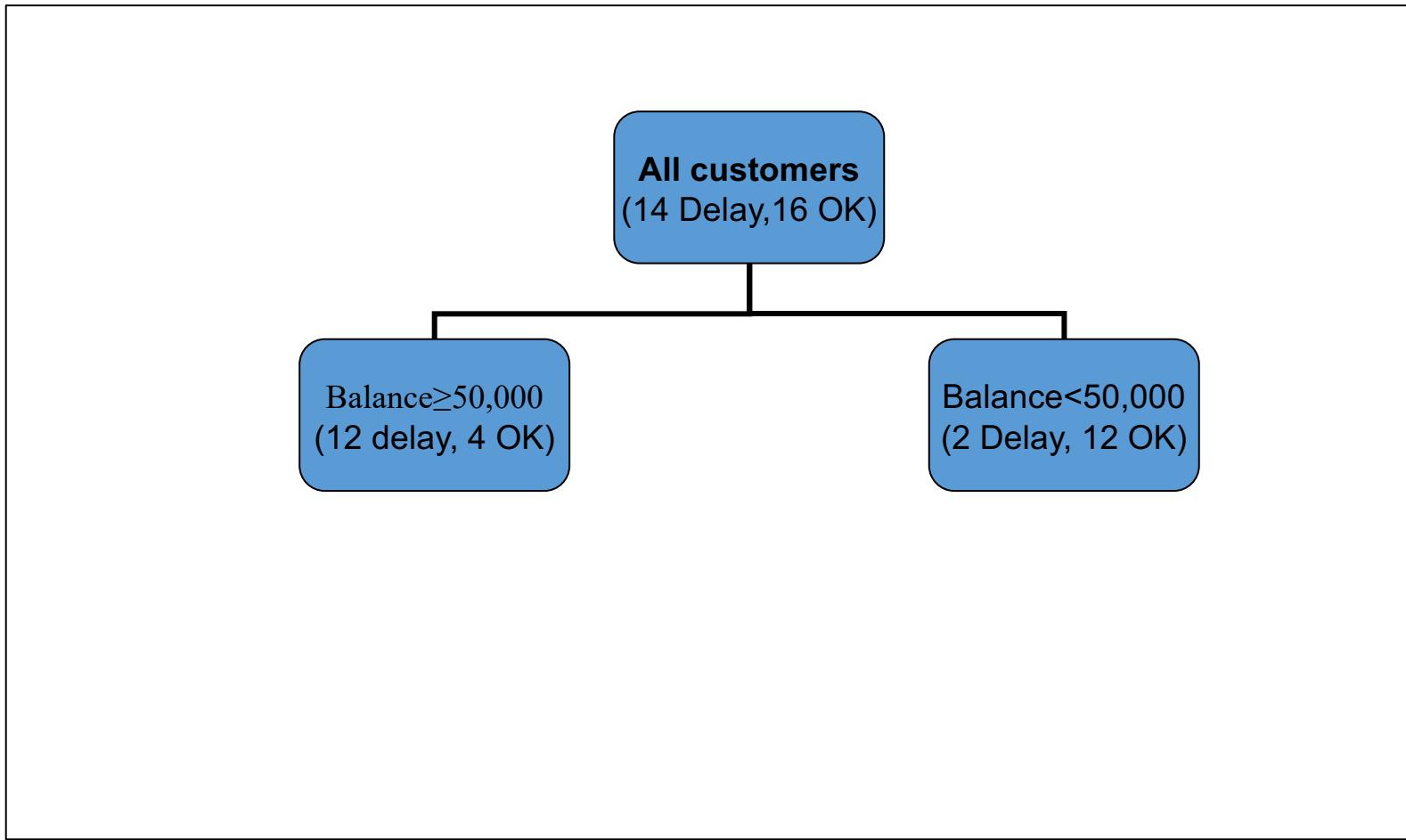


All customers
(14 Delay, 16 OK)

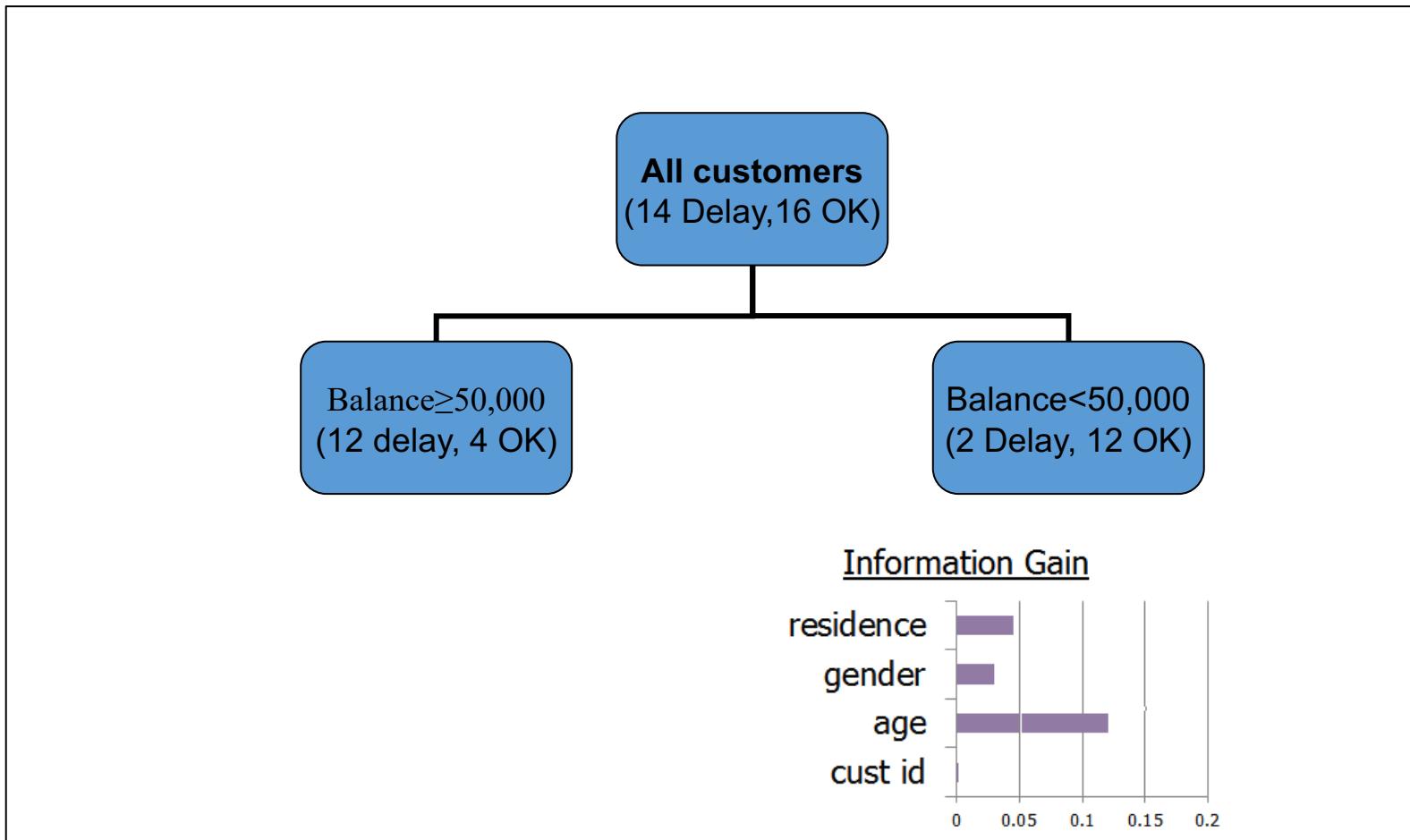
Tree Structure



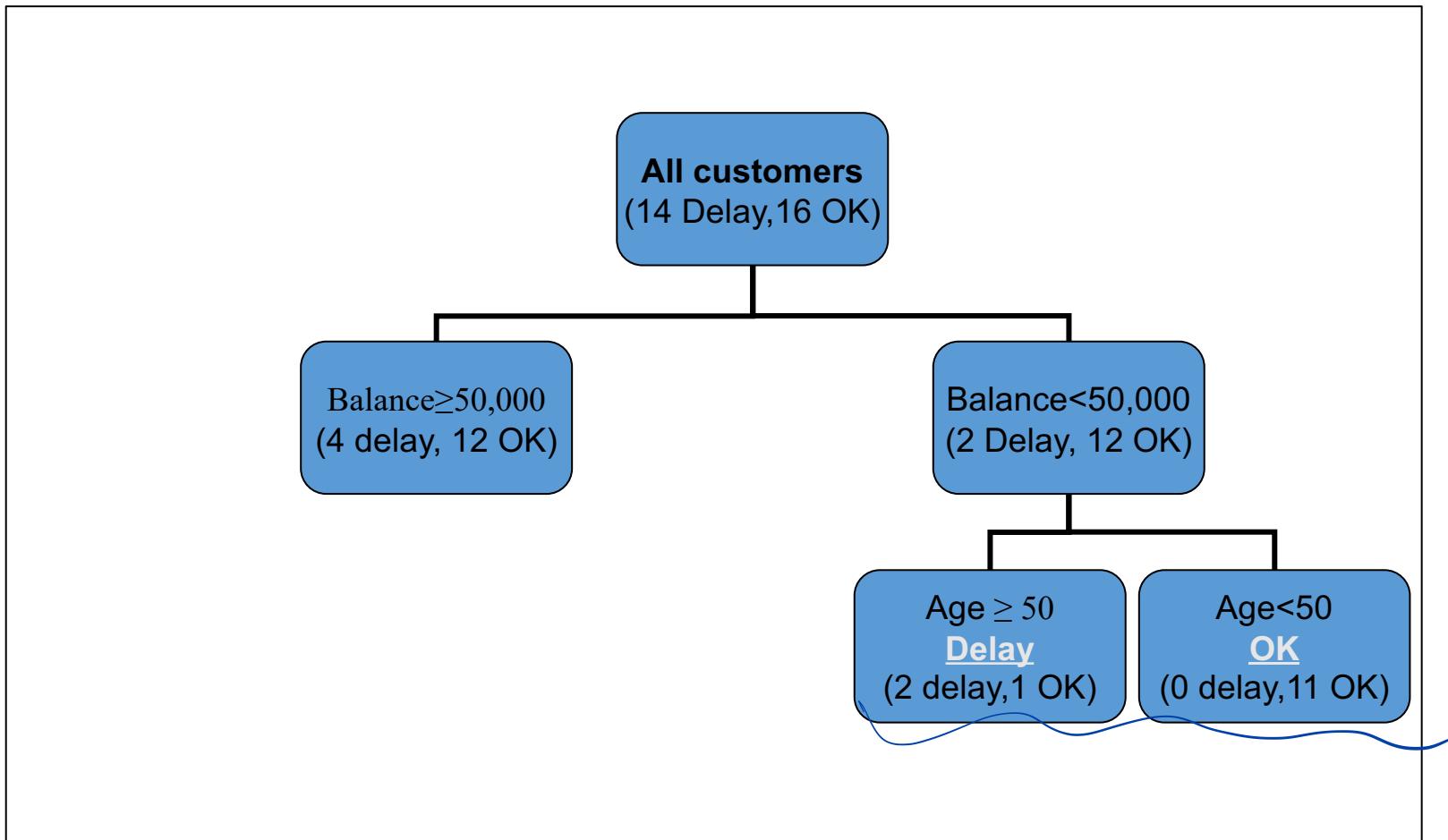
Tree Structure



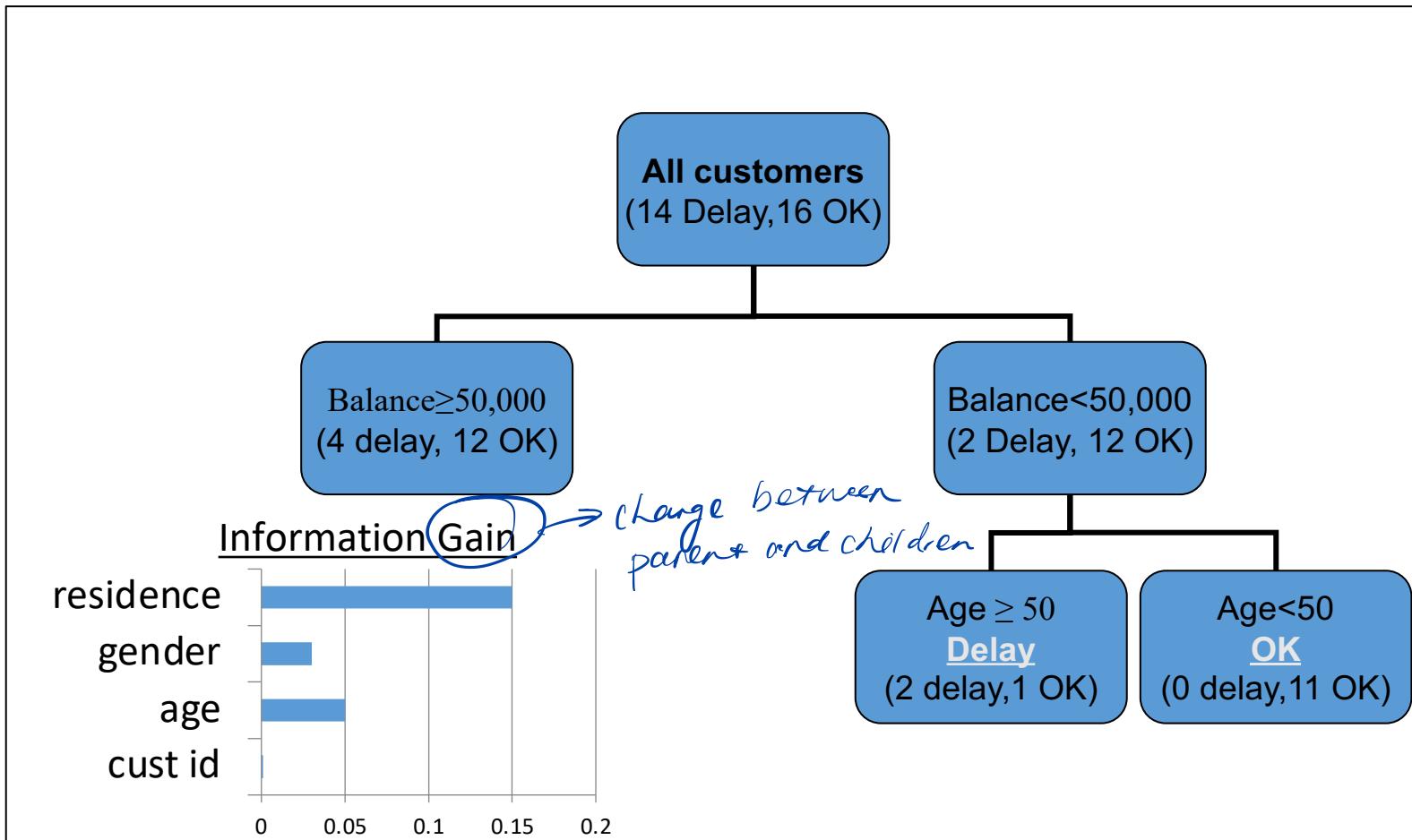
Tree Structure



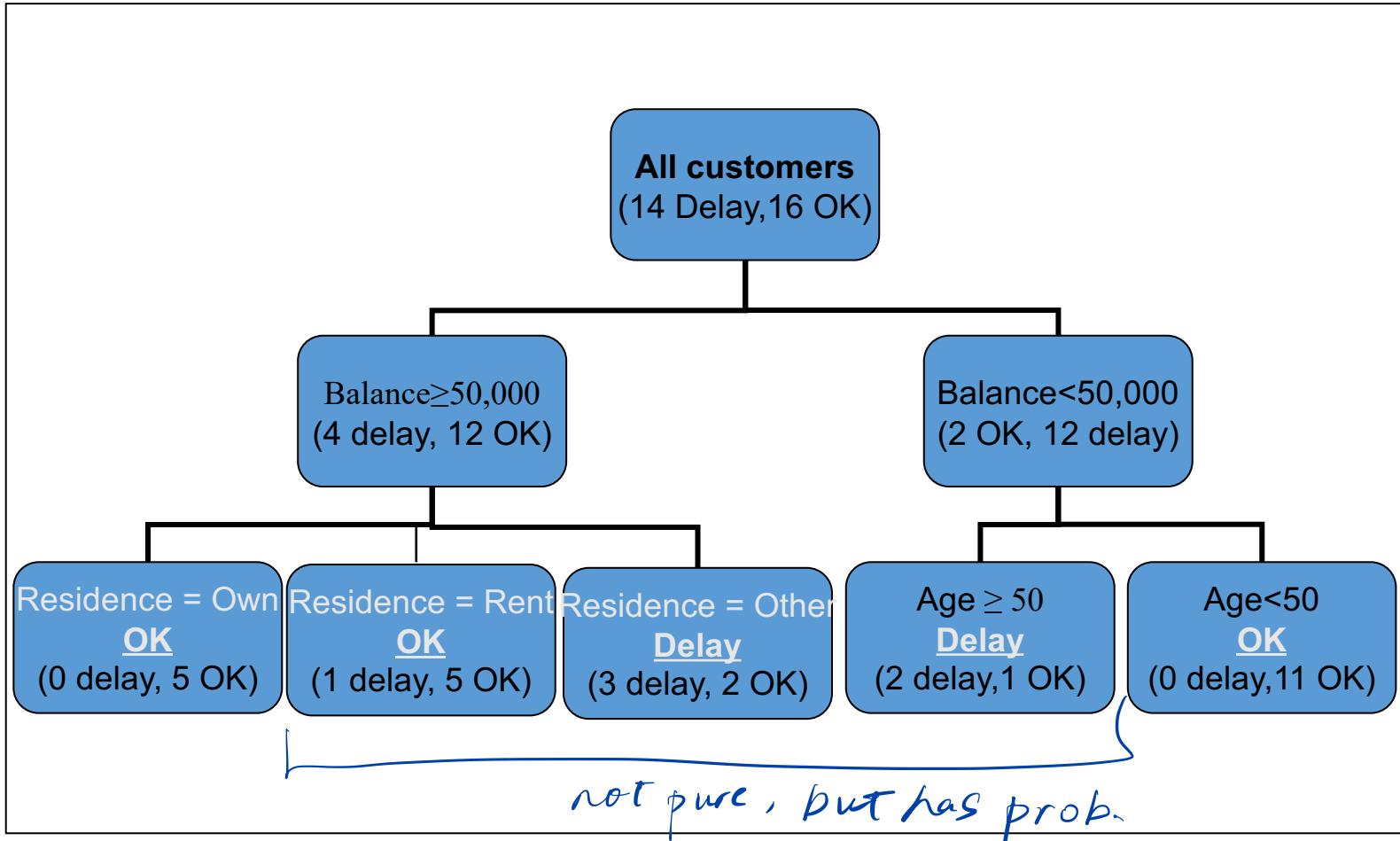
Tree Structure



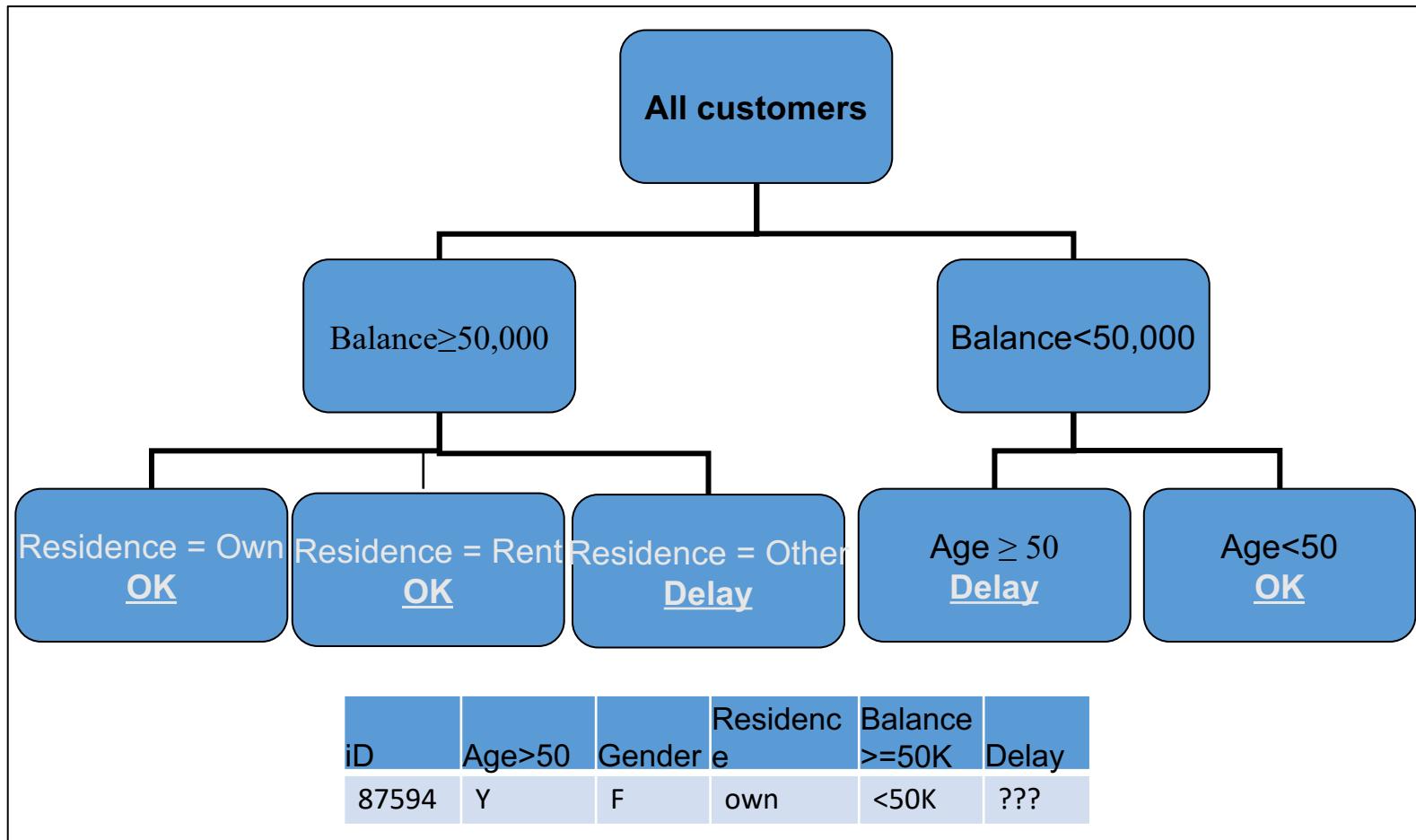
Tree Structure



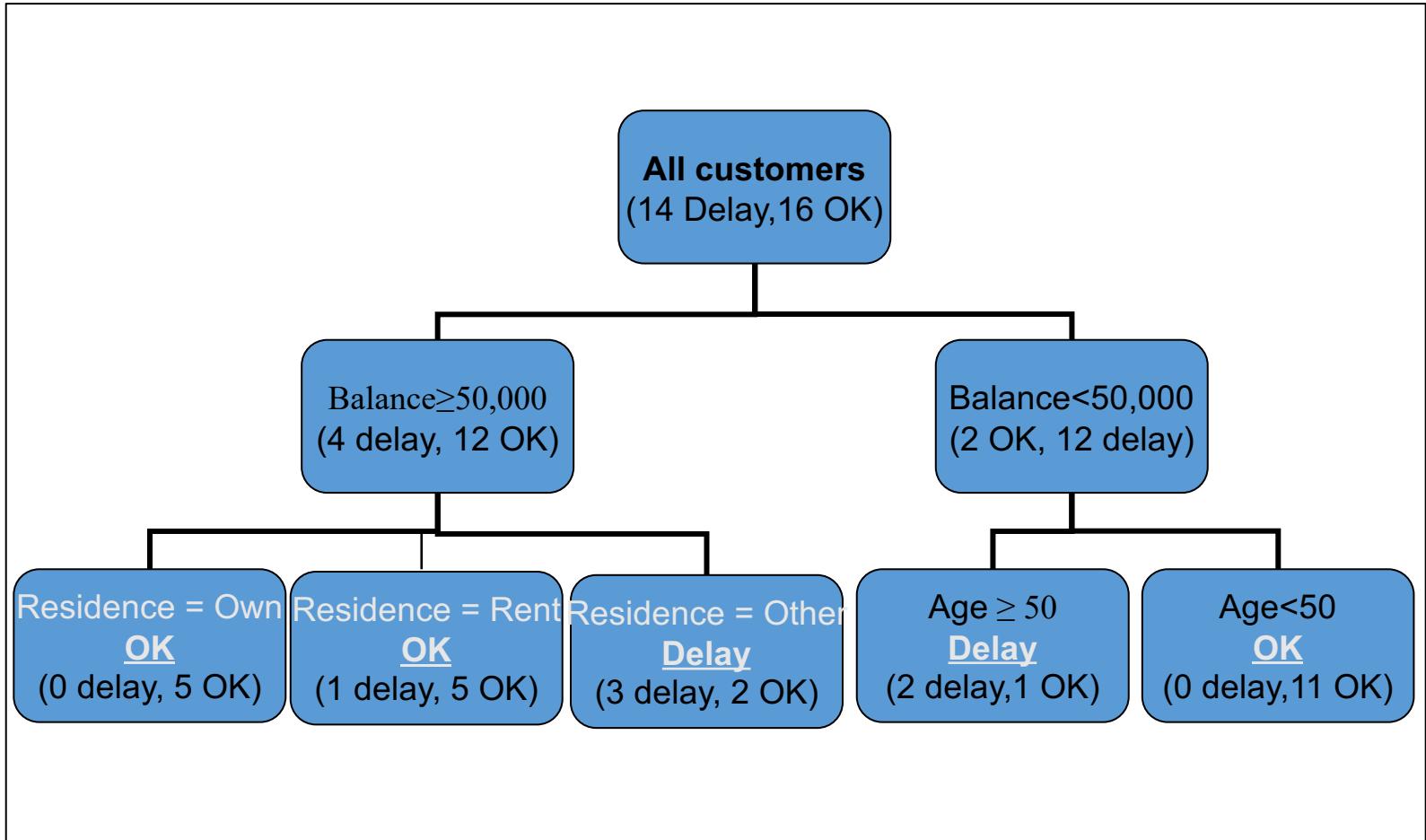
Tree Structure



Tree Structure



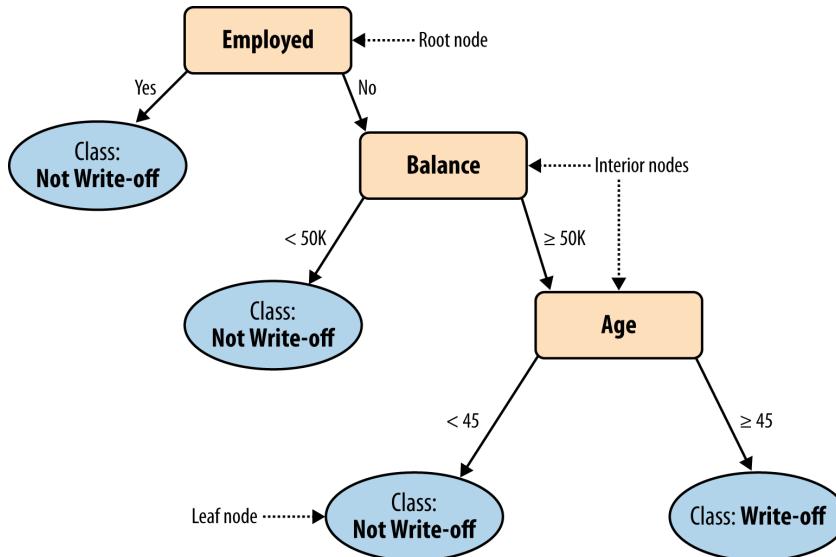
Open Issues



Trees as Sets of Rules

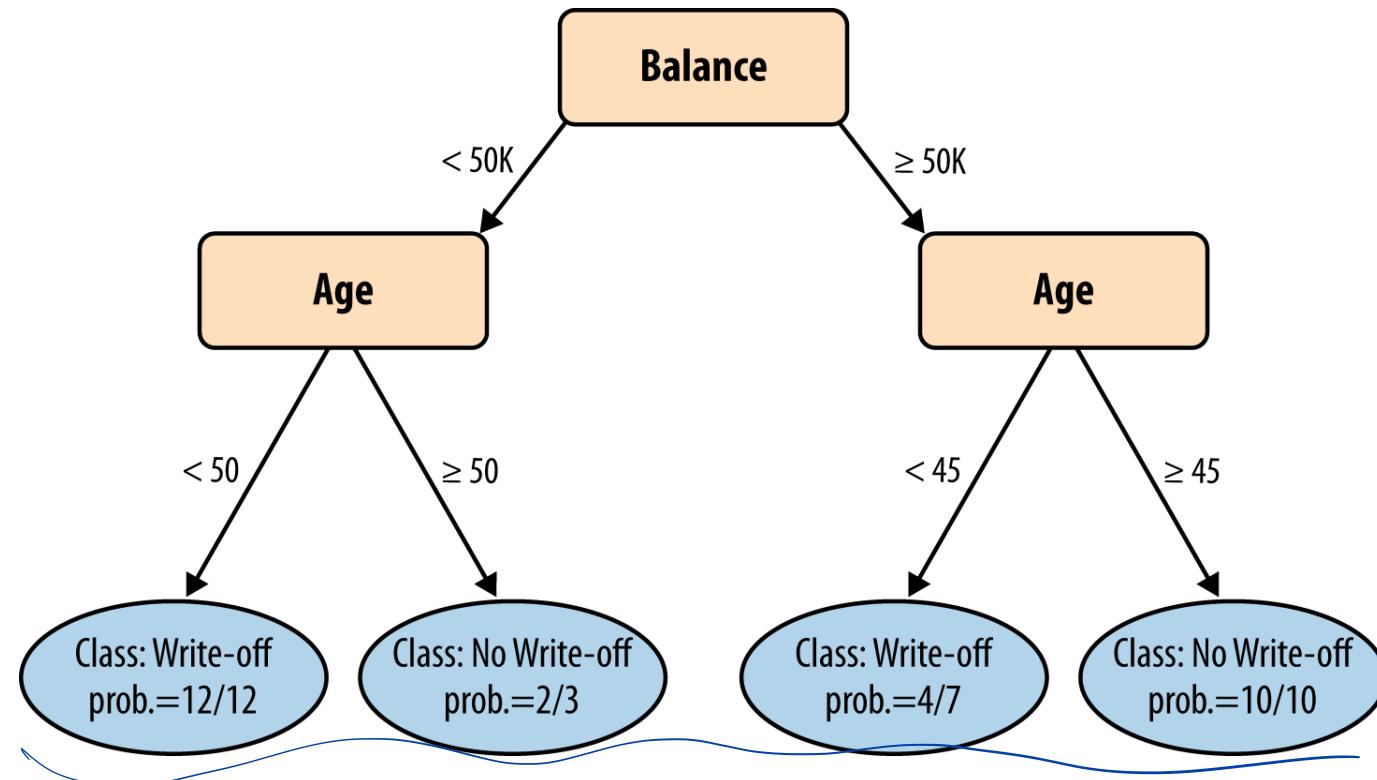
- The classification tree is equivalent to set of rules
- Each rule consists of the attribute tests along the path connected with **AND**

Trees as Sets of Rules



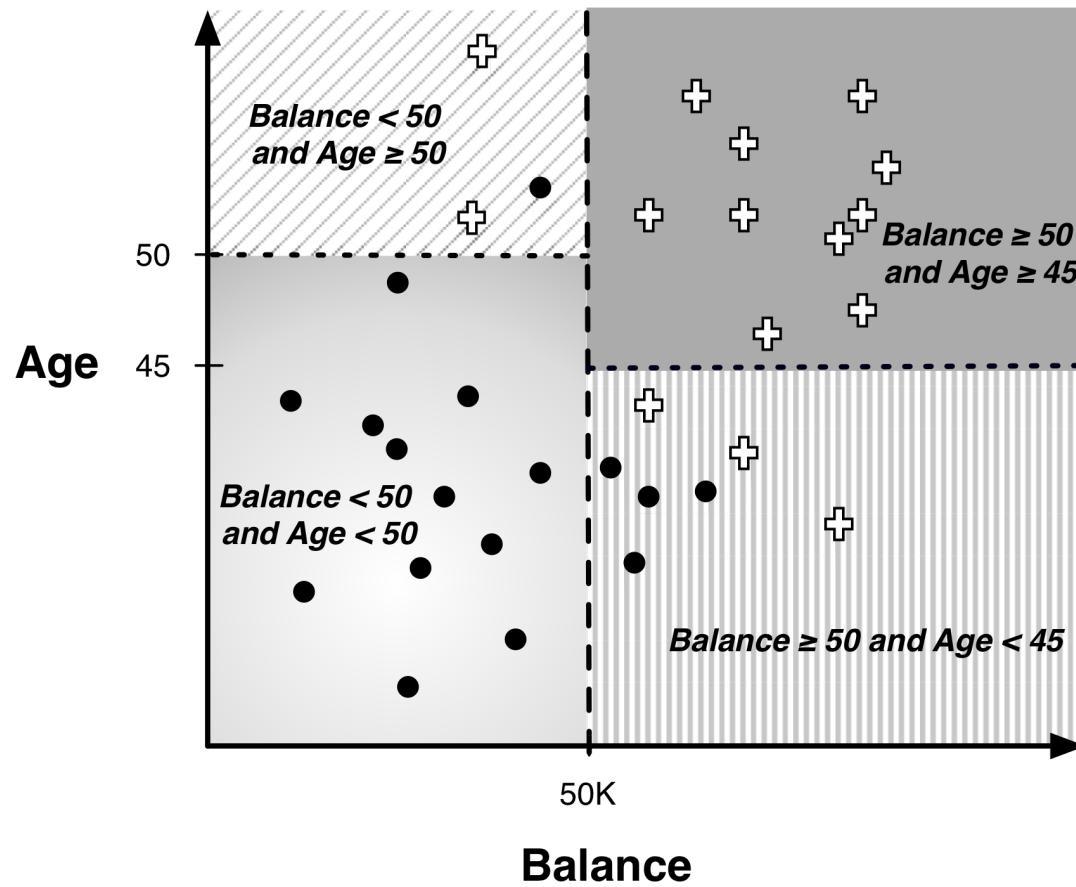
- IF (Employed = Yes) THEN Class=No Write-off
- IF (Employed = No) AND (Balance < 50k) THEN Class=No Write-off
- IF (Employed = No) AND (Balance ≥ 50k) AND (Age < 45) THEN Class=No Write-off
- IF (Employed = No) AND (Balance ≥ 50k) AND (Age ≥ 45) THEN Class=Write-off

Visualizing Segmentations



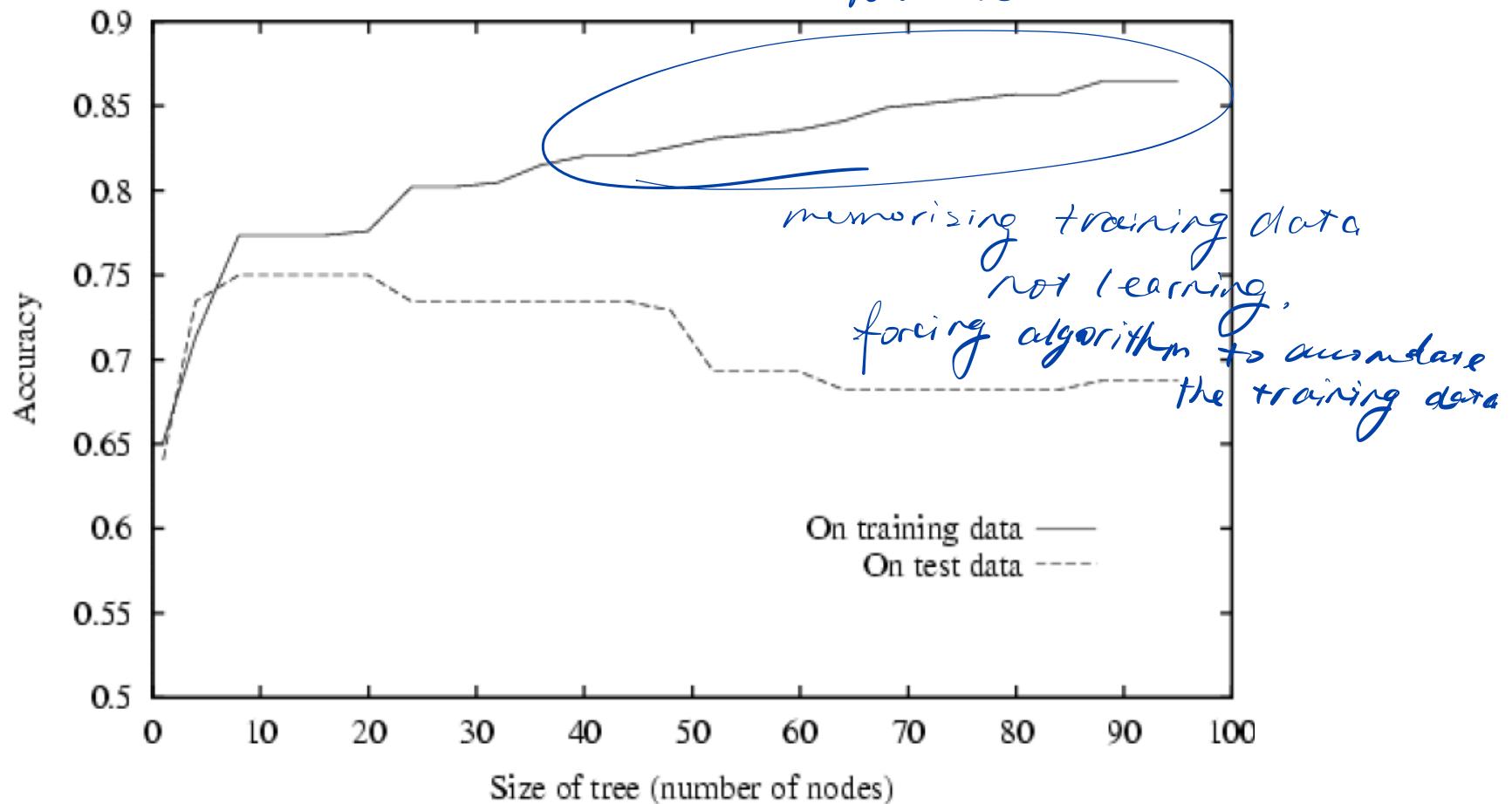
Visualizing Segmentations

如何多个会成为高维

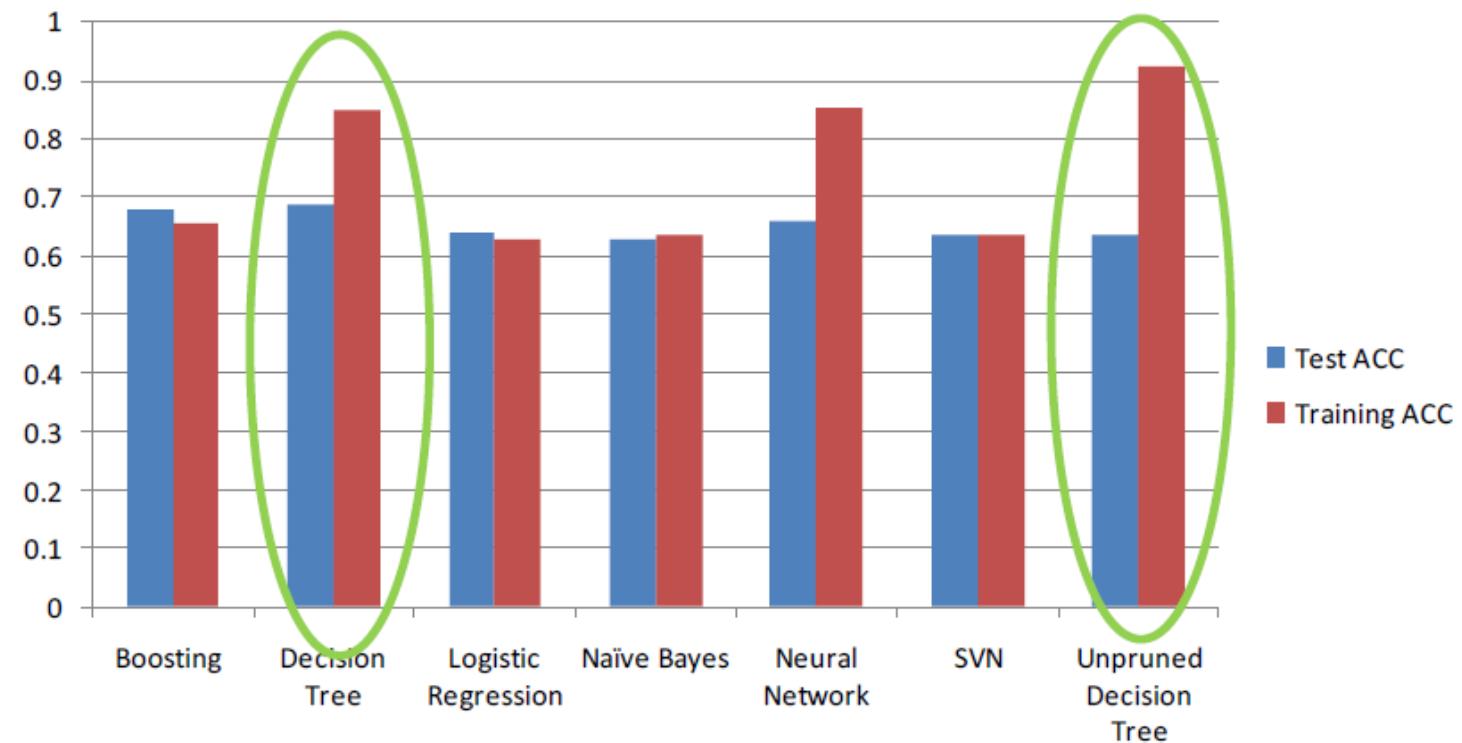


Tree Complexity and Over-fitting

constraint : attributes, sample size,
node 節点数



Trees on Churn



Pruning

- Pruning simplifies a decision tree to prevent over-fitting to noise in the data
- **Post-pruning:**
 - takes a fully-grown decision tree and discards unreliable parts
- **Pre-pruning:**
 - stops growing a branch when information becomes unreliable
- Post-pruning preferred in practice

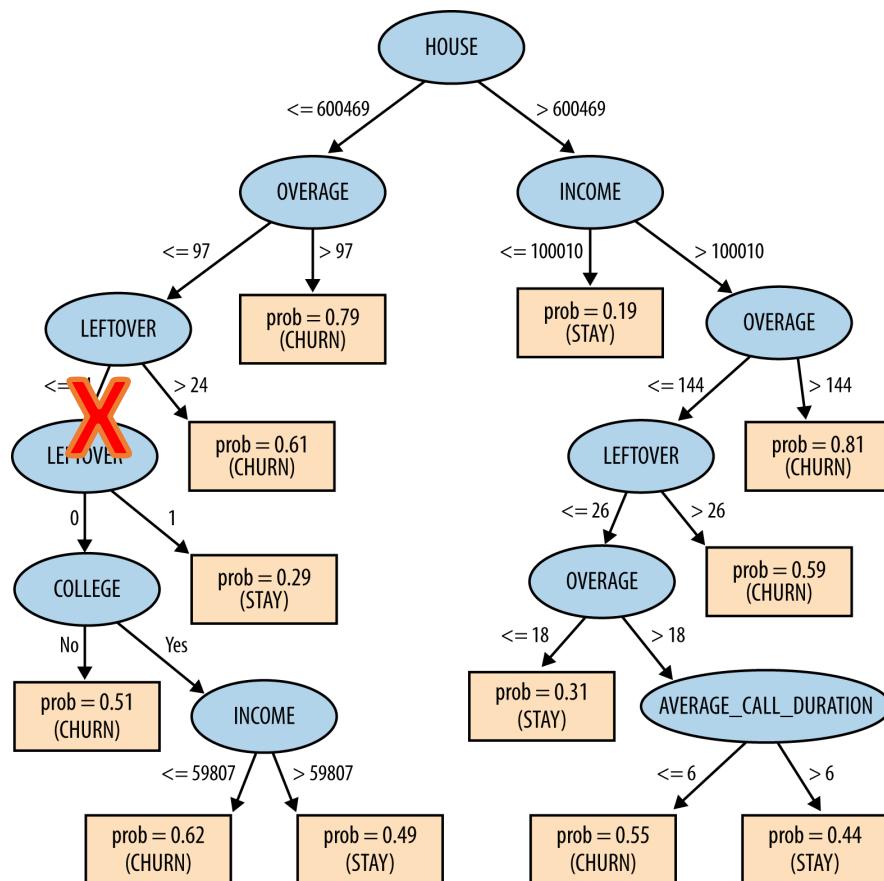
Decision Tree Pruning Methods

- Validation set – withhold a subset (~1/3) of training data to use for pruning
 - Note: you should randomize the order of training examples

Reduced-Error Pruning

- Classify examples in validation set – some might be errors
- For each node:
 - Sum the errors over entire subtree
 - Calculate error on same example if converted to a leaf with majority class label
- Prune node with highest reduction in error
- Repeat until error no longer reduced

MegaTelCo: Predicting Churn with Tree Induction



From Classification Trees to Probability Estimation Trees

- **Frequency-based estimate**
 - **Basic assumption:** Each member of a segment corresponding to a tree leaf has the same probability to belong in the corresponding class
 - If a leaf contains n positive instances and m negative instances (binary classification), the probability of any new instance being positive may be estimated as $\frac{n}{n+m}$
- Prone to **over-fitting..**

Laplace Correction

- $p(c) = \frac{n+1}{n+m+2}$,
 - where n is the number of examples in the leaf belonging to class c , and m is the number of examples not belonging to class c

The many faces of classification: Classification / Probability Estimation / Ranking

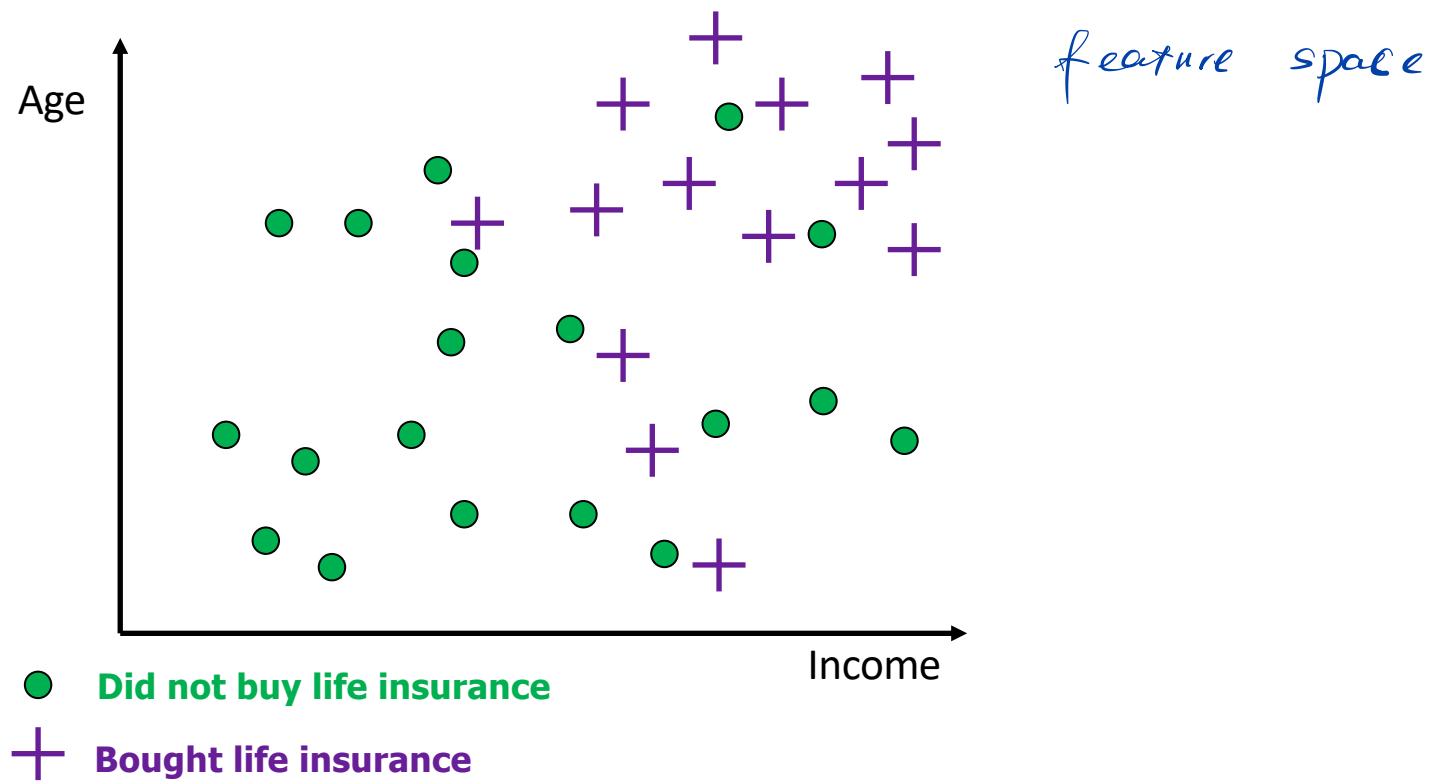
- Classification Problem
 - Most general case: The target takes on discrete values that are **NOT** ordered
 - Most common: binary classification where the target is either 0 or 1
- 3 Different Solutions to Classification
 - Classifier model: Model predicts the same set of **discrete value** as the data had
 - In binary case:
 - **Ranking:** Model predicts a **score** where a higher score indicates that the model think the example to be more likely to be in one class
 - **Probability estimation:** Model predicts a **score between 0 and 1** that is meant to be the probability of being in that class

The many faces of classification: Classification / Probability Estimation / Ranking

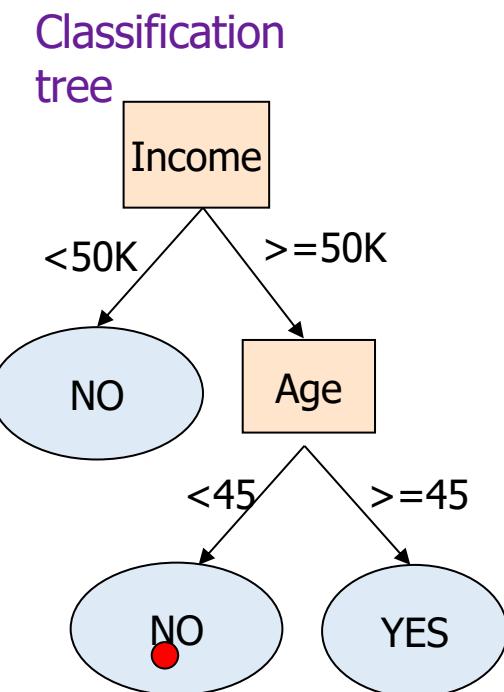
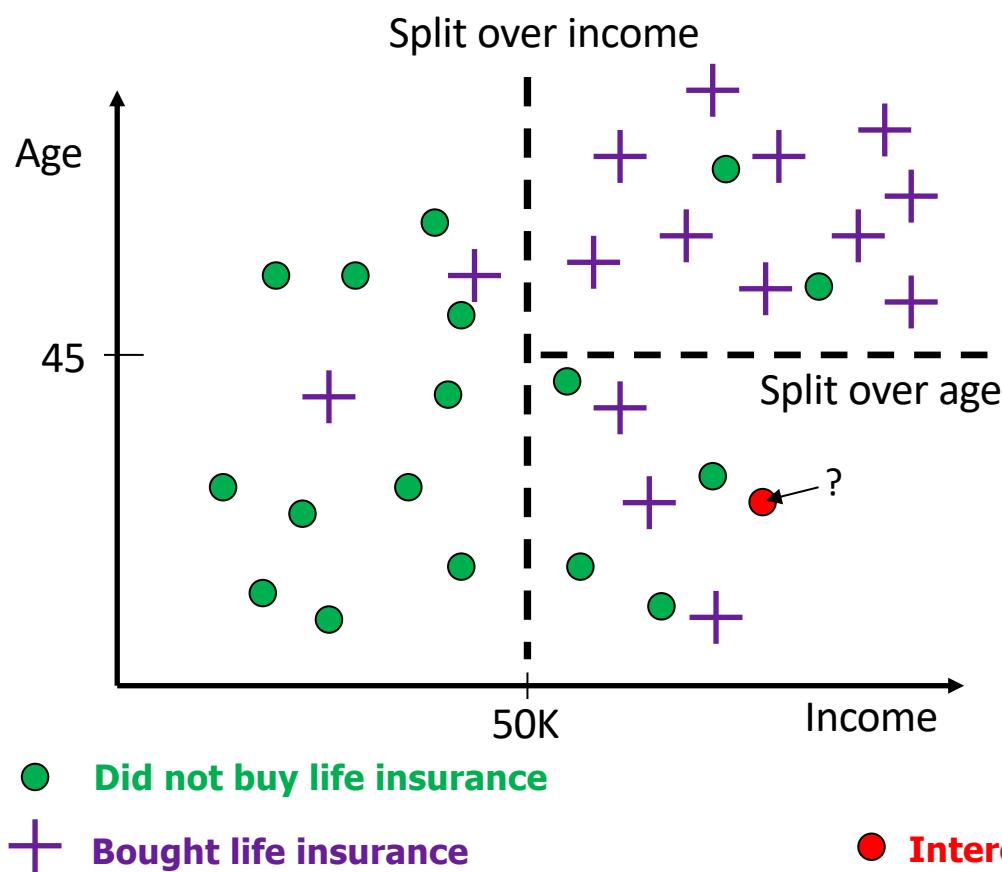


- Ranking:
 - business context determines the number of actions (“how far down the list”)
 - cost/benefit is constant, unknown, or difficult to calculate
- Probability:
 - you can always rank / classify if you have probabilities!
 - cost/benefit is not constant across examples and known relatively precisely

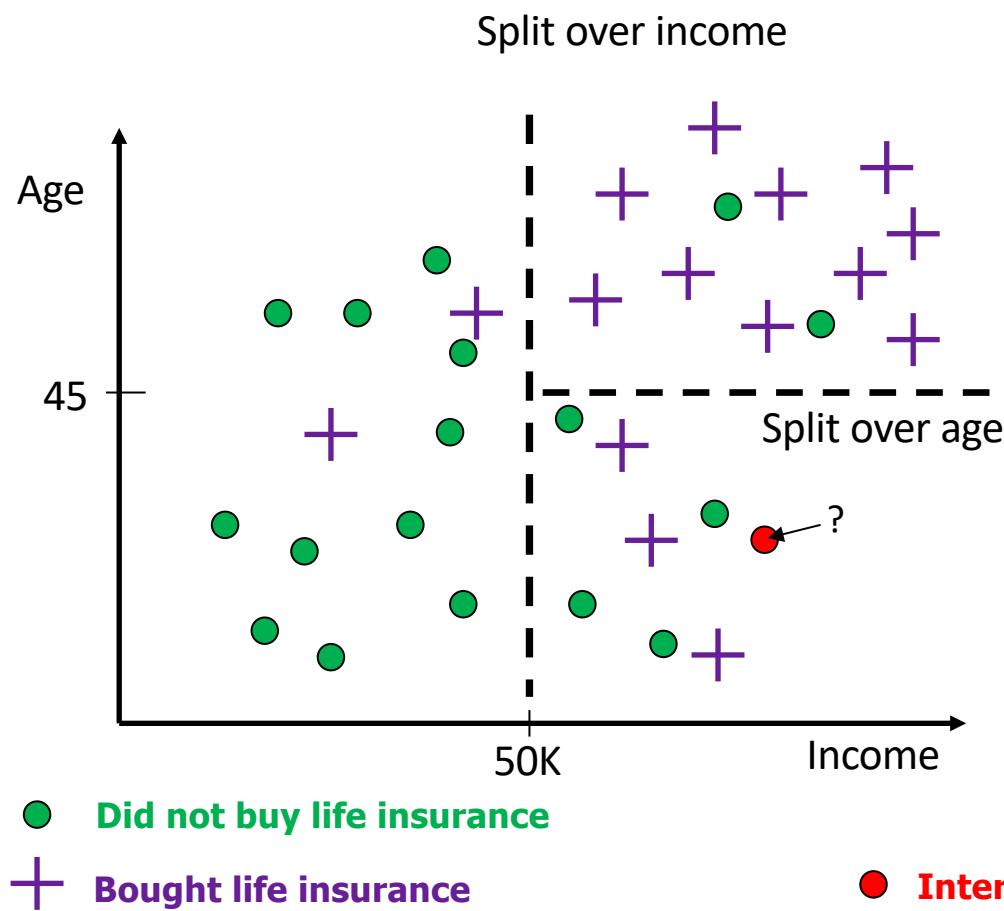
Example



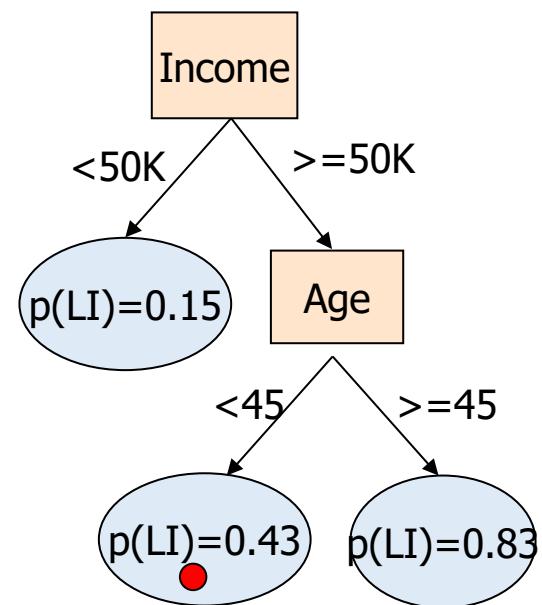
Example



Example



Classification tree

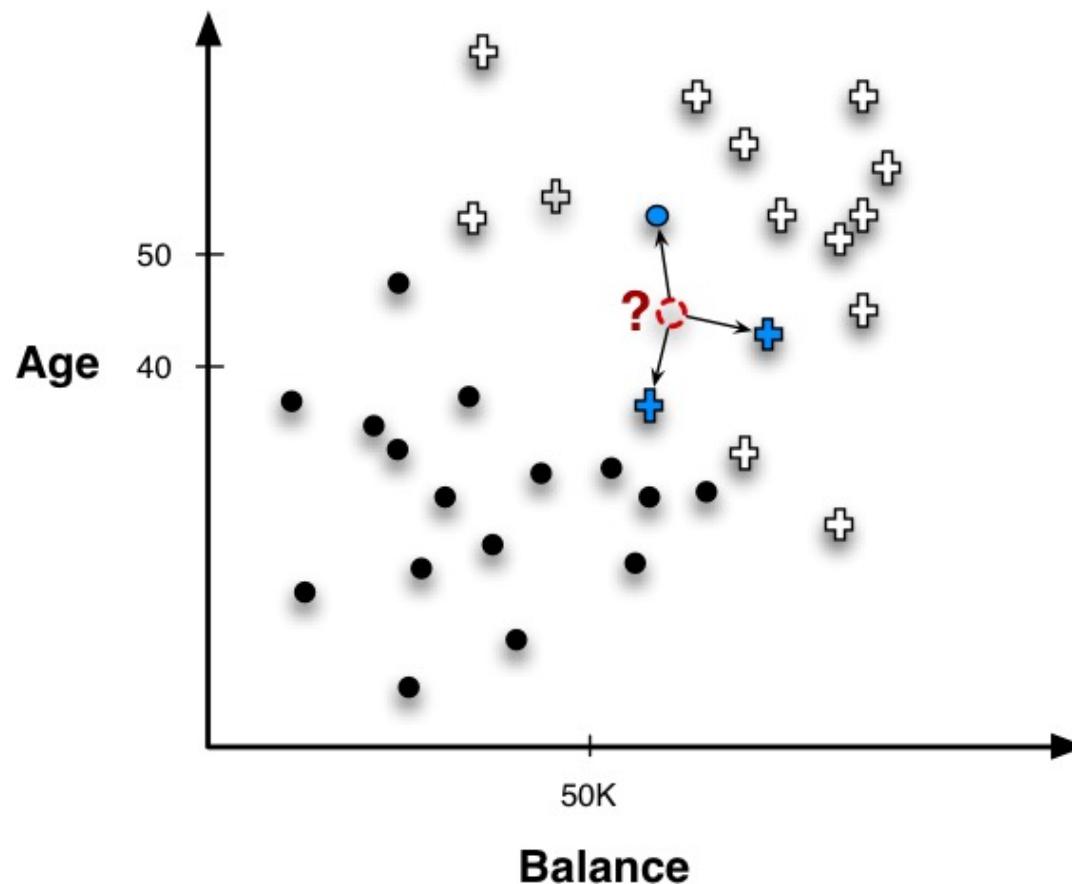


● Interested in LI? = 3/7

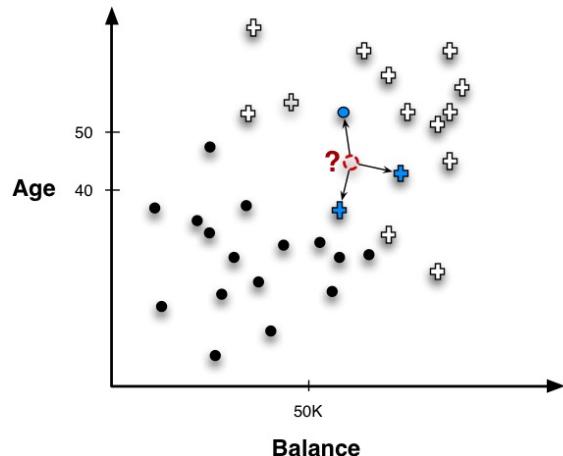
Class Outline

- Decision Trees
- Nearest Neighbor Classification
- Algorithmic Analysis
- Ensemble Methods
- Evaluation Metrics

Nearest neighbors for predictive modeling



Nearest neighbor in classification



- ▶ Majority vote
- ▶ How many neighbors?
- ▶ Also consider the distance to vary the influence of the neighbors

Customer	Age	Income (1000s)	Cards	Response (target)	Distance from David
David	37	50	2	?	
John	35	35	3	Yes	$\sqrt{(35 - 37)^2 + (35 - 50)^2 + (3 - 2)^2} = 15.16$
Rachael	22	50	2	No	$\sqrt{(22 - 37)^2 + (50 - 50)^2 + (2 - 2)^2} = 15$
Ruth	63	200	1	No	$\sqrt{(63 - 37)^2 + (200 - 50)^2 + (1 - 2)^2} = 152.23$
Jefferson	59	170	1	No	$\sqrt{(59 - 37)^2 + (170 - 50)^2 + (1 - 2)^2} = 122$
Norah	25	40	4	Yes	$\sqrt{(25 - 37)^2 + (40 - 50)^2 + (4 - 2)^2} = 15.74$

Name	Distance	Similarity	Weight	Contribution	Class
Rachael	15.0		0.004444	0.344	No
John	15.2		0.004348	0.336	Yes
Norah	15.7		0.004032	0.312	Yes
Jefferson	122.0		0.000067	0.005	No
Ruth	152.2		0.000043	0.003	No

Nearest neighbor classification

Given a labeled training set $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})$

Example: the MNIST dataset of handwritten digits.

A 10x10 grid of handwritten digits, likely from the MNIST dataset. The digits are rendered in black on a white background. They vary in style and orientation, some appearing more upright than others. The grid is composed of ten rows and ten columns of digits.

1	4	1	6	1	1	9	1	5	4
8	6	6	3	5	9	7	2	0	2
0	1	3	0	8	4	1	1	1	5
3	1	1	0	6	4	1	1	0	3
6	6	8	9	1	2	0	7	4	7
6	0	2	0	1	8	7	3	0	1
8	4	0	1	0	9	7	0	7	5
5	5	1	0	7	5	5	1	8	2
4	3	1	7	8	7	5	4	1	6
5	5	1	8	2	5	5	1	0	8

To classify a new instance x :

Find its nearest neighbor amongst the $x^{(i)}$, Return $y^{(i)}$

The data space

We need to choose a distance function.



Each image is 28×28 grayscale.
One option: Treat images as 784-dimensional vectors, and use Euclidean (ℓ_2) distance:

$$\|x - x'\| = \sqrt{\sum_{i=1}^{784} (x_i - x'_i)^2}.$$

Summary:

- Data space $X = \mathbb{R}^{784}$ with ℓ_2 distance
- Label space $Y = \{0, 1, \dots, 9\}$

Performance on MNIST

Training set of 60,000 points.

- What is the error rate on training points?

Performance on MNIST

Training set of 60,000 points.

- What is the error rate on training points? **Zero.**
In general, **training error** is an overly optimistic predictor of future performance.

Performance on MNIST

Training set of 60,000 points.

- What is the error rate on training points? **Zero**.
In general, **training error** is an overly optimistic predictor of future performance.
- A better gauge: separate test set of 10,000 points.
Test error = fraction of test points incorrectly classified.

Performance on MNIST

Training set of 60,000 points.

- What is the error rate on training points? **Zero**.
In general, **training error** is an overly optimistic predictor of future performance.
- A better gauge: separate test set of 10,000 points.
Test error = fraction of test points incorrectly classified.
- What test error would we expect for a random classifier?

Performance on MNIST

Training set of 60,000 points.

- What is the error rate on training points? **Zero**.
In general, **training error** is an overly optimistic predictor of future performance.
- A better gauge: separate test set of 10,000 points.
Test error = fraction of test points incorrectly classified.
- What test error would we expect for a random classifier? **90%**

Performance on MNIST

Training set of 60,000 points.

- What is the error rate on training points? **Zero**.
In general, **training error** is an overly optimistic predictor of future performance.
- A better gauge: separate test set of 10,000 points.
Test error = fraction of test points incorrectly classified.
- What test error would we expect for a random classifier? **90%**
- Test error of nearest neighbor: **3.09%**.

Performance on MNIST

Training set of 60,000 points.

- What is the error rate on training points? **Zero**.
In general, **training error** is an overly optimistic predictor of future performance.
- A better gauge: separate test set of 10,000 points.
Test error = fraction of test points incorrectly classified.
- What test error would we expect for a random classifier? **90%**
- Test error of nearest neighbor: **3.09%**.

Examples of errors:



Performance on MNIST

Training set of 60,000 points.

- What is the error rate on training points? **Zero**.
In general, **training error** is an overly optimistic predictor of future performance.
- A better gauge: separate test set of 10,000 points.
Test error = fraction of test points incorrectly classified.
- What test error would we expect for a random classifier? **90%**
- Test error of nearest neighbor: **3.09%**.

Examples of errors:



Ideas for improvement: (1) k -NN (2) better distance function.

K -nearest neighbor classification

Classify a point using the labels of its k nearest neighbors among the training points.

K -nearest neighbor classification

Classify a point using the labels of its k nearest neighbors among the training points.

MNIST:	k	1	3	5	7	9	11
	Test error (%)	3.09	2.94	3.13	3.10	3.43	3.34

K -nearest neighbor classification

Classify a point using the labels of its k nearest neighbors among the training points.

MNIST:	k	1	3	5	7	9	11
	Test error (%)	3.09	2.94	3.13	3.10	3.43	3.34

How to choose k in general?

K -nearest neighbor classification

Classify a point using the labels of its k nearest neighbors among the training points.

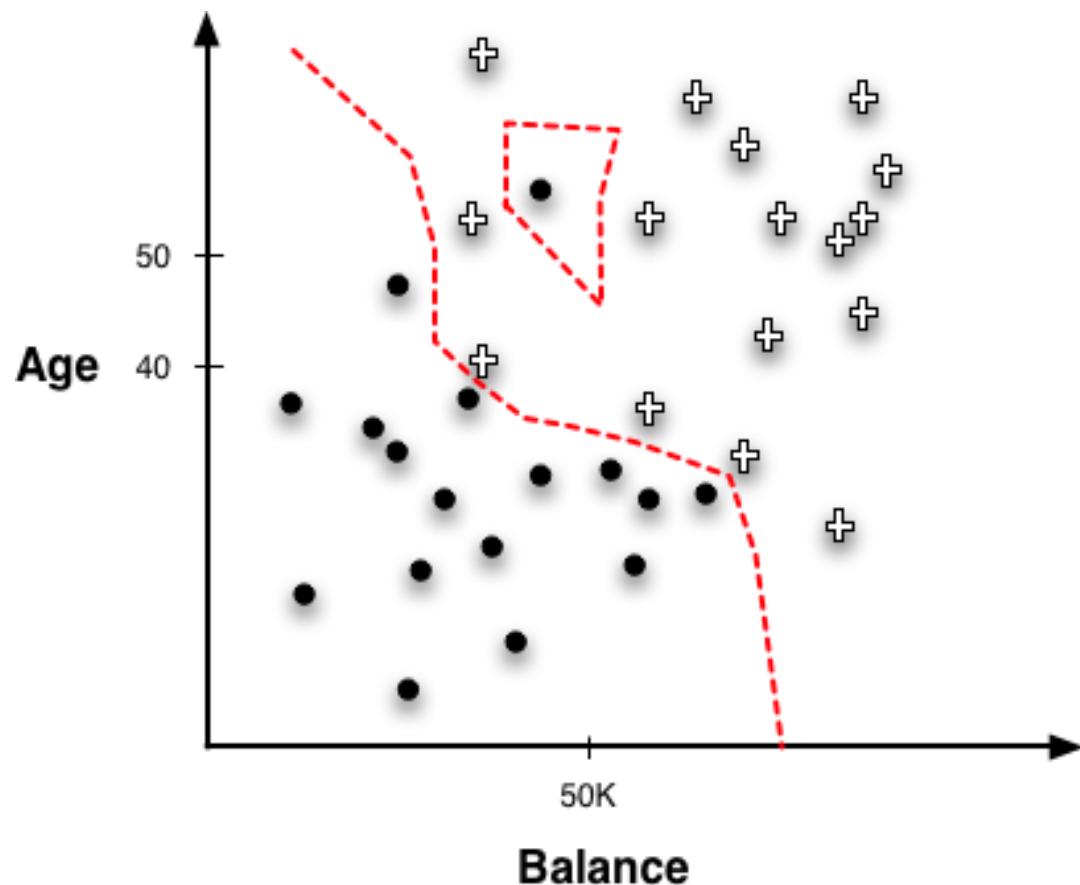
MNIST:	k	1	3	5	7	9	11
	Test error (%)	3.09	2.94	3.13	3.10	3.43	3.34

How to choose k in general?

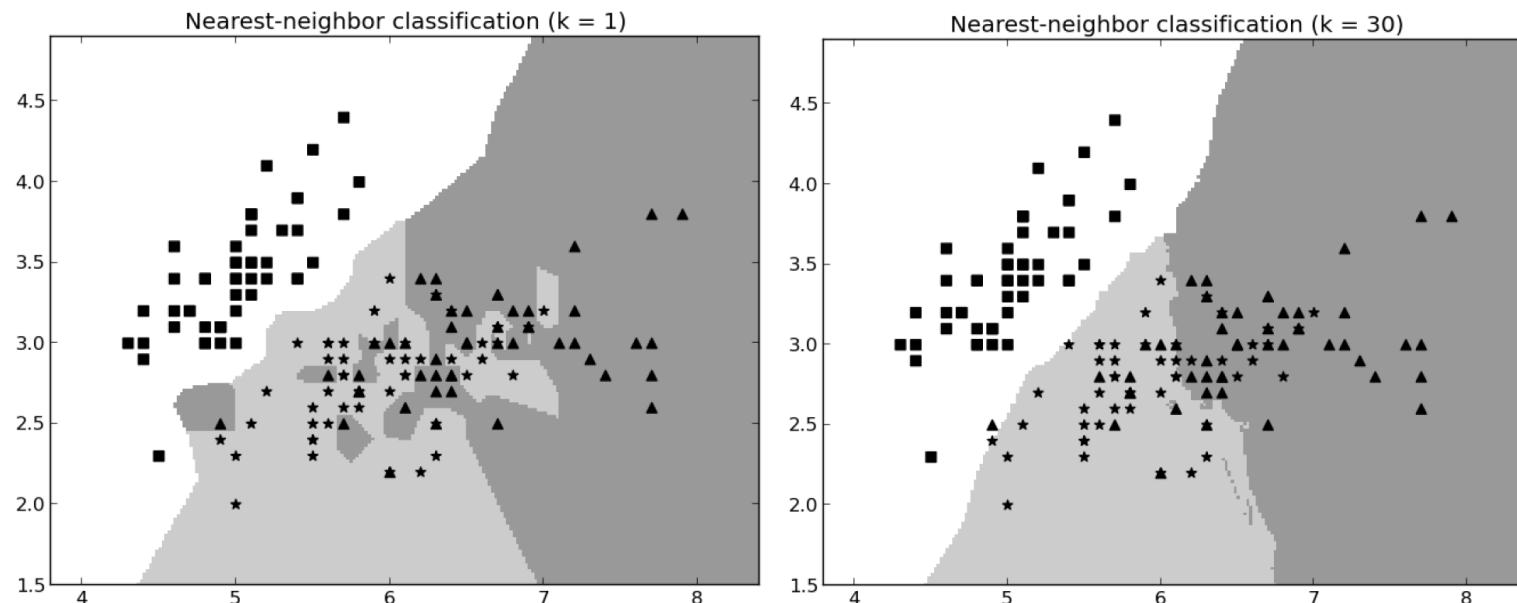
- Let $S \in Z^n$ be the training set, where $Z = X \times Y$ is the space of *labeled* points.
- Select a **Hold-out set** V – also called the **Validation set**
- Train the k-NN on $(S - V)$ for some value of k
- Check the performance of this classifier on V
- Repeat for a different value of k
- Pick the best k (one with minimum test error)

The above procedure can also be performed using **Leave-one-out cross validation** in which we choose one training point at a time as the validation set

Geometric Interpretation, Over-fitting, and Complexity



Nearest neighbor in classification



- ▶ Nearest neighbor classifiers follow very specific boundaries
- ▶ 1-NN strongly tends to overfit (k is a complexity parameter!)
- ▶ Use cross-validation or nested holdout testing

Better distance functions

Let x be an image. Consider an image x' that is just like x , but is either:

- shifted one pixel to the right, or
- Rotated slightly.

Then $\|x - x'\|$ could easily be quite large.

Better distance functions

Let x be an image. Consider an image x' that is just like x , but is either:

- shifted one pixel to the right, or
- Rotated slightly.

Then $\|x - x'\|$ could easily be quite large.

It makes sense to choose distance measures that are invariant under:

- Small translations and rotations. E.g. *tangent distance*.
- A broader family of natural deformations. E.g. *shape context*

Better distance functions

Let x be an image. Consider an image x' that is just like x , but is either:

- shifted one pixel to the right, or
- Rotated slightly.

Then $\|x - x'\|$ could easily be quite large.

It makes sense to choose distance measures that are invariant under:

- Small translations and rotations. E.g. *tangent distance*.
- A broader family of natural deformations. E.g. *shape context*

Test error rates:	ℓ_2	tangent distance	shape context
	3.09	1.10	0.63

Better distance functions

Tangent Distance

- Is an invariant distance measure
- Especially effective for OCR
- Small transformations of certain image objects does not affect class membership

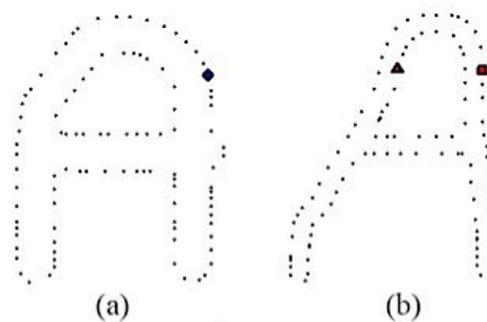
$$s(\text{[image of a handwritten digit]}, \alpha) = \begin{matrix} \text{[image of a handwritten digit]} \\ \alpha = -2 \end{matrix} \quad \begin{matrix} \text{[image of a handwritten digit]} \\ \alpha = -1 \end{matrix} \quad \begin{matrix} \text{[image of a handwritten digit]} \\ \alpha = 0 \end{matrix} \quad \begin{matrix} \text{[image of a handwritten digit]} \\ \alpha = 1 \end{matrix} \quad \begin{matrix} \text{[image of a handwritten digit]} \\ \alpha = 2 \end{matrix}$$

Refer: <http://yann.lecun.com/exdb/publis/pdf/simard-00.pdf> for more details

Better distance functions

Shape Context

- Is a feature descriptor in object recognition
- A way of describing shapes that allows for measuring shape similarity and the recovering of point correspondences
- pick n points on the contours of a shape



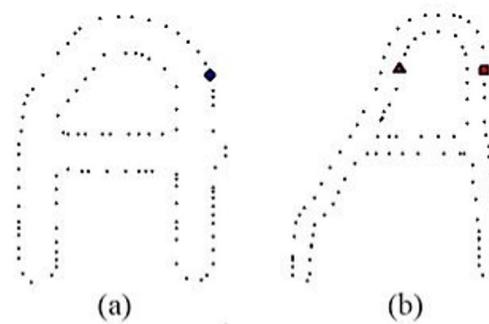
References:

<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.441.6897&rep=rep1&type=pdf>
https://en.wikipedia.org/wiki/Shape_context

Better distance functions

Shape Context

- Is a feature descriptor in object recognition
- A way of describing shapes that allows for measuring shape similarity and the recovering of point correspondences
- pick n points on the contours of a shape



Are there other families of distance functions that are often useful?

ℓ_p norms

How can we measure the length of a vector in \mathbb{R}^m ?

Usual choice is the *Euclidean norm*:

$$\|x\|_2 = \sqrt{\sum_{i=1}^m x_i^2}.$$

ℓ_p norms

How can we measure the length of a vector in \mathbb{R}^m ?

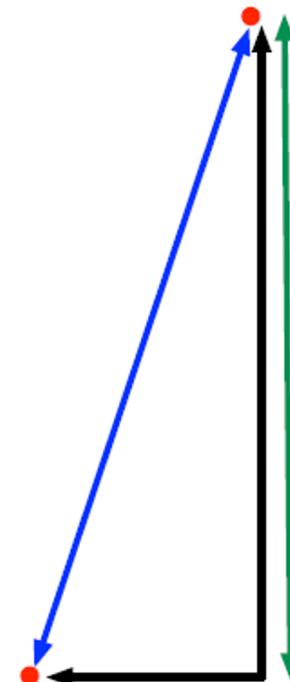
Usual choice is the *Euclidean norm*:

$$\|x\|_2 = \sqrt{\sum_{i=1}^m x_i^2}.$$

Generalization: For $p \geq 1$, the ℓ_p norm is

$$\|x\|_p = \left(\sum_{i=1}^m |x_i|^p \right)^{1/p}$$

- $p = 2$: Euclidean norm
- ℓ_1 norm: $\|x\|_1 = \sum_{i=1}^m |x_i|$
- ℓ_∞ norm: $\|x\|_\infty = \max_i |x_i|$



Quick quiz

Suppose data lie in \mathbb{R}^p

- 1 What is the l_2 norm of the all-ones vector?

Quick quiz

Suppose data lie in \mathbb{R}^p

- ① What is the l_2 norm of the all-ones vector?
- ② Suppose $\|x\|_1 = 1$.

What is the maximum value of the l_2 norm?

Quick quiz: Answers

Suppose data lie in \mathbb{R}^p

- ① What is the l_2 norm of the all-ones vector? $1/(p)^{1/2}$
- ② Suppose $\|x\|_1 = 1$.

What is the maximum value of the l_2 norm? **One.**

Metric spaces

A more general notion is a *metric space*.

Metric spaces

A more general notion is a *metric space*.

Let X be the space in which data lie. A distance function $d : X \times X \rightarrow \mathbb{R}$ is a *metric* if it satisfies these properties:

- $d(x, y) \geq 0$ (non-negativity)
- $d(x, y) = 0$ if and only if $x = y$
- $d(x, y) = d(y, x)$ (symmetry)
- $d(x, z) \leq d(x, y) + d(y, z)$ (triangle inequality)

Metric spaces

A more general notion is a *metric space*.

Let X be the space in which data lie. A distance function $d : X \times X \rightarrow \mathbb{R}$ is a *metric* if it satisfies these properties:

- $d(x, y) \geq 0$ (non-negativity)
- $d(x, y) = 0$ if and only if $x = y$
- $d(x, y) = d(y, x)$ (symmetry)
- $d(x, z) \leq d(x, y) + d(y, z)$ (triangle inequality)

For instance:

- $\mathcal{X} = \mathbb{R}^m$ and $d(x, y) = \|x - y\|_p$
- $\mathcal{X} = \{\text{strings over some alphabet}\}$ and $d = \text{edit distance}$.

Class Outline

- Decision Trees
- Nearest Neighbor Classification
- Algorithmic Analysis
- Ensemble Methods
- Evaluation Metrics

Sensitivity to noise

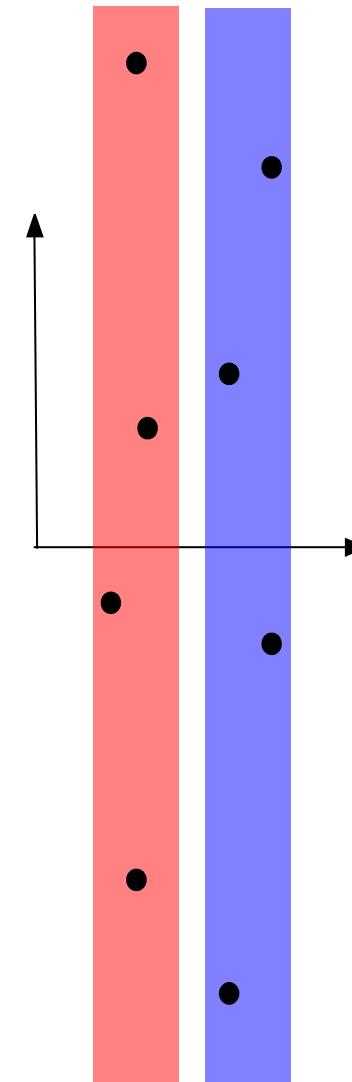
Adding a single sufficiently noisy feature can wreak havoc with nearest neighbor classifiers.

Sensitivity to noise

Adding a single sufficiently noisy feature can wreak havoc with nearest neighbor classifiers.



versus



How to do feature selection/reweighting
for nearest neighbor?

Class Outline

- Decision Trees
- Nearest Neighbor Classification
- Algorithmic Analysis
- **Ensemble Methods**
- Evaluation Metrics

Bagging

Given a data set S of n labeled points:

- For $t = 1$ to T :
 - Choose n' points randomly, with replacement, from S .
 - Fit a classifier h_t to these points.

(For instance: $T = 100$ or 1000 , $n' = n$.)

Final predictor: majority vote of h_1, \dots, h_T

The resampling and averaging reduces overfitting.

Random forests

Rather like bagging decision trees, but with an additional source of randomization.

Given a data set S of n labeled points:

- For $t = 1$ to T :
 - Choose n' points randomly, with replacement, from S .
 - Fit a decision tree h_t to these points. When building the tree, at each node restrict the split direction to be one of k features at random. chosen at random.

(For instance: $k = \sqrt{p}$ when data is p -dimensional.)

Final predictor: majority vote of h_1, \dots, h_T .

Class Outline

- Decision Trees
- Nearest Neighbor Classification
- Algorithmic Analysis
- Ensemble Methods
- Evaluation Metrics

Evaluation Metrics

- Up to now: measure a model's performance by some simple metric
 - classifier error rate, accuracy, ...
- Simple example: accuracy

$$accuracy = \frac{\text{Number of correct decisions made}}{\text{Total number of decisions made}}$$

- Classification accuracy is popular, but usually **too simplistic** for applications of data mining to real business problems
- **Decompose** and count the different types of correct and incorrect decisions made by a classifier

Unequal costs and benefits

- How much do we care about the different **errors** and correct decisions?
 - Classification accuracy makes no distinction between **false positive** and **false negative** errors
 - In real-world applications, different kinds of errors lead to different consequences!
- Examples for medical diagnosis:
 - a patient has cancer (although he does not)
→ **false positive error**, expensive, but not life threatening
 - a patient has cancer, but she is told that she has not
→ **false negative error**, more serious
- Errors should be counted separately
 - Estimate cost or benefit of each decision

Confusion Matrix

- A **confusion matrix** for a problem involving n classes
 - is an $n \times n$ matrix with the columns labeled with actual classes and the rows labeled with predicted classes

		Predicted Classes	
		p	n
True Values	1	True Positives (TP)	False Negative (FN)
	0	False Positives (FP)	True Negatives (TN)

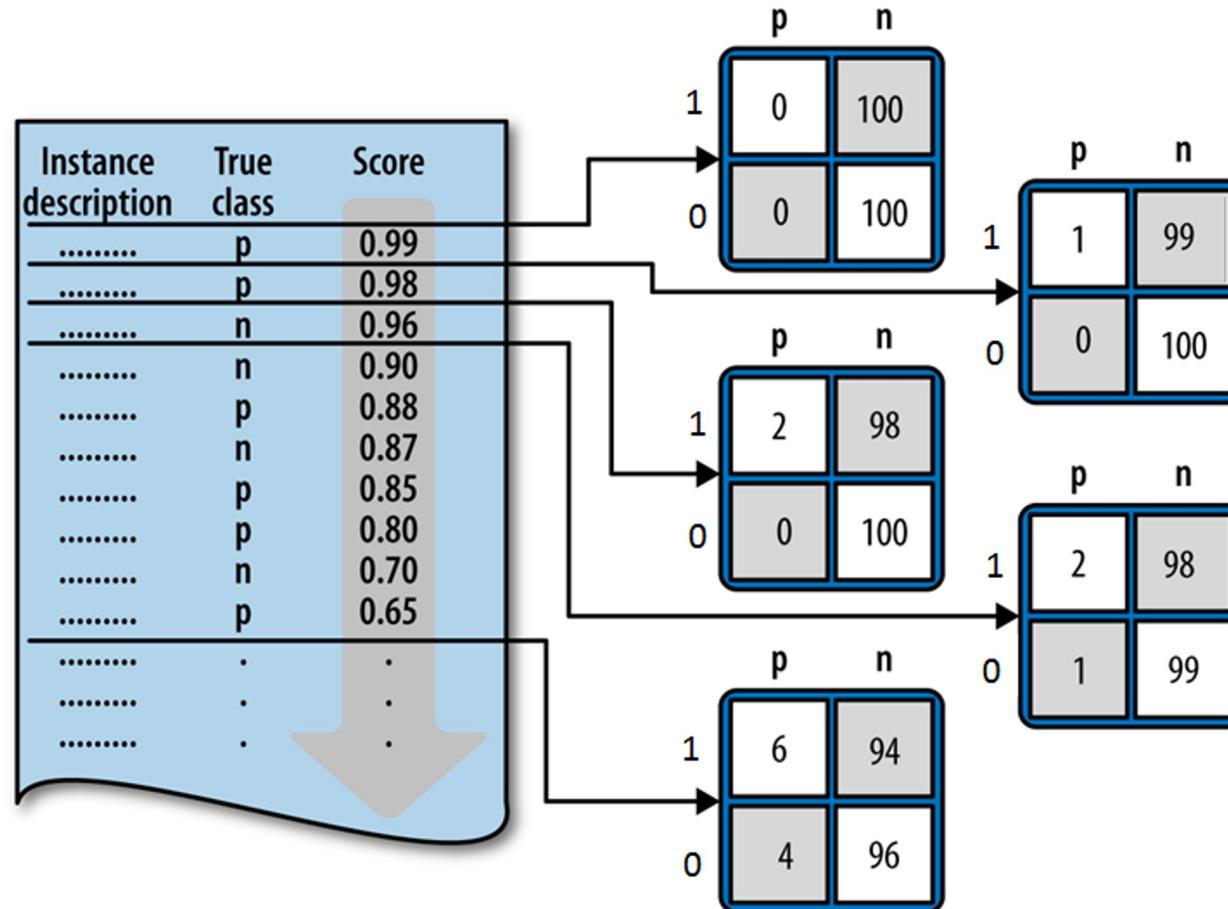
- Each example in a test set has an **actual class label** and the **class predicted** by the classifier
- The confusion matrix separates out the decisions made by the classifier
 - actual/true classes: **1** (Positive Label), **0** (Negative Label)
 - predicted classes: **p**(ositive), **n**(egative)
 - The main diagonal contains the count of correct decisions

Other Evaluation Metrics

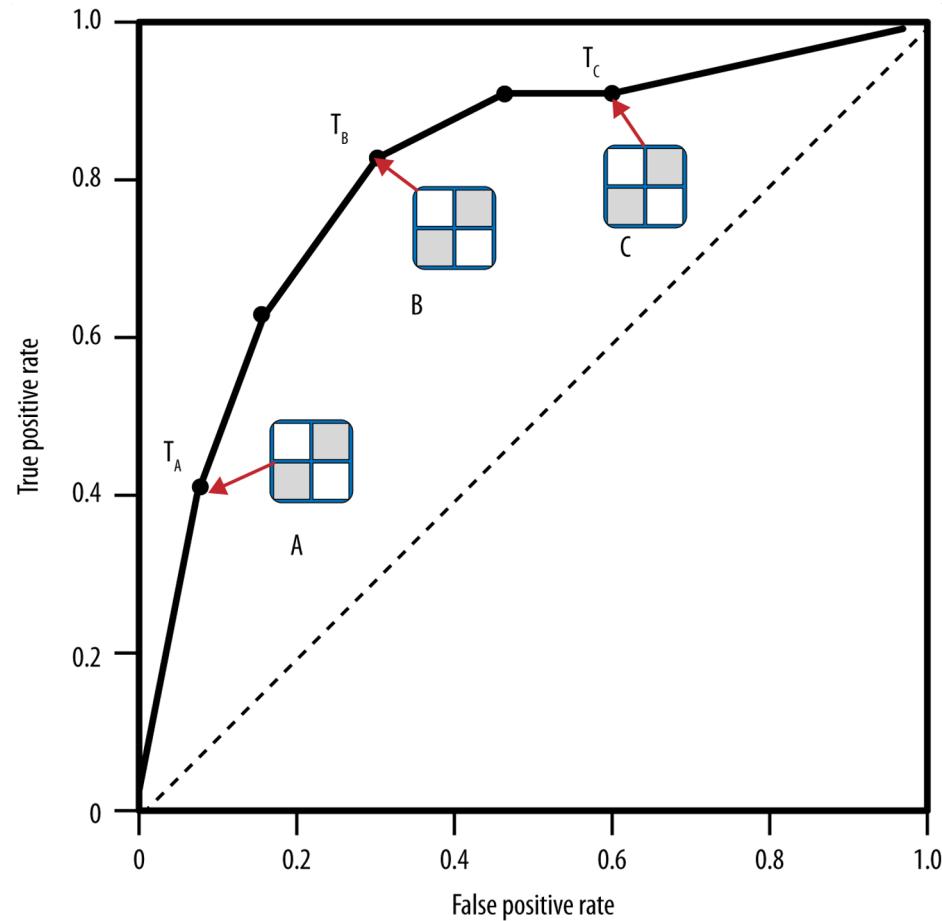
- Based on the entries of the confusion matrix, we can describe various evaluation metrics
 - True positive rate (Recall): $\frac{TP}{TP+FN}$
 - False negative rate: $\frac{FN}{TP+FN}$
 - Precision (accuracy over the cases predicted to be positive): $\frac{TP}{TP+FP}$
 - F-measure (harmonic mean): $2 \cdot \frac{precision \cdot recall}{precision + recall}$
 - Specificity: $\frac{TN}{TN+FP}$
 - Sensitivity: $\frac{TP}{TP+FN}$
 - Accuracy (count of correct decisions): $\frac{TP+TN}{P+N}$
 - False Positive Rate = $1 - \text{Sensitivity} = \frac{FN}{TP+FN}$

ROC Curve

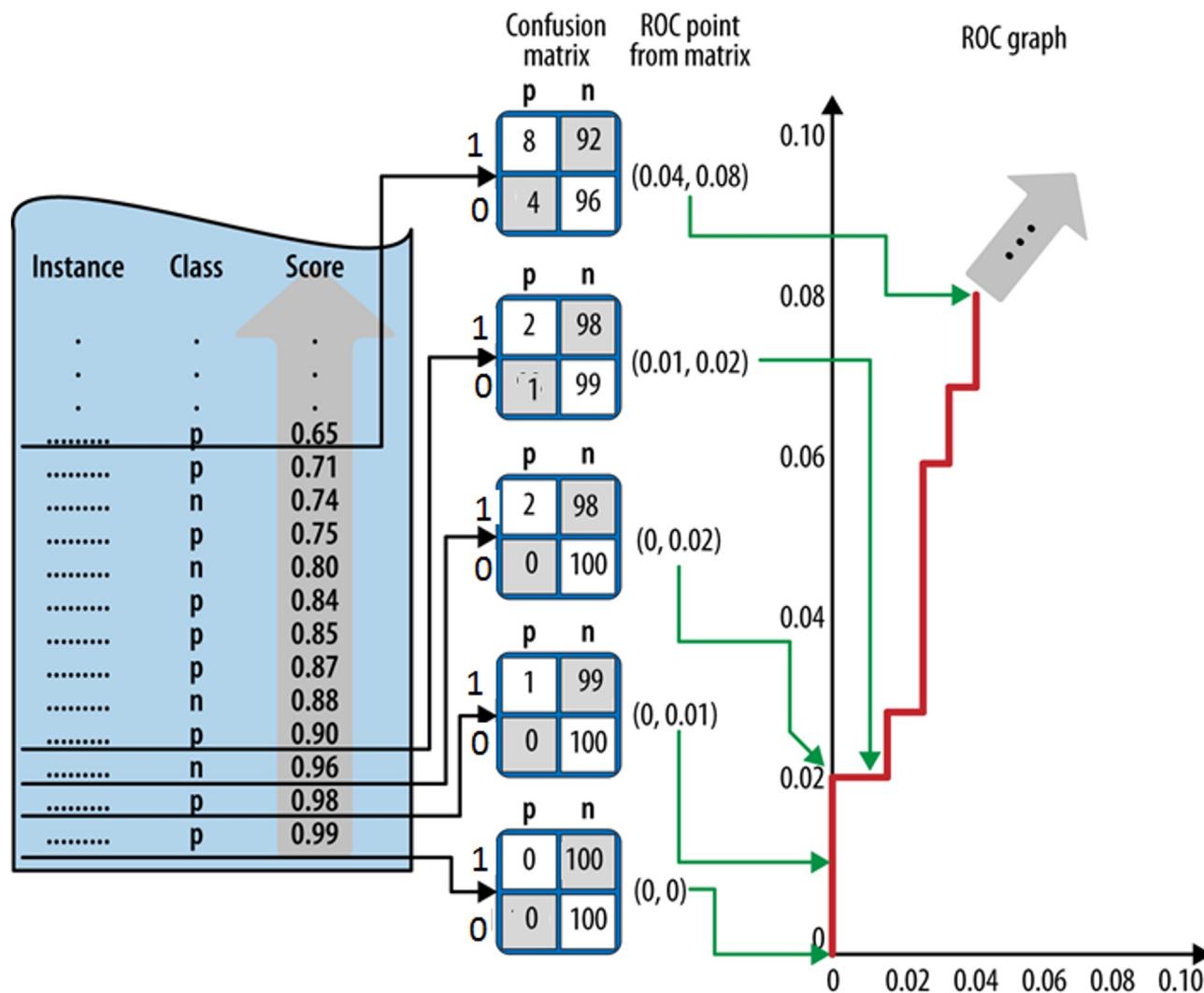
Ranking instead of classifying



ROC graphs and curves



ROC graphs and curves



Getting ROC curve: Algorithm

- Sort the test set by the model predictions
- Start with cutoff = max (prediction)
- Decrease cutoff, after each step count the number of true positives TP (positives with prediction above the cutoff) and false positives FP (negatives above the cutoff)
- Calculate TP rate (TP/P) and FP (FP/N) rate
- Plot current number of TP/P as a function of current FP/N

Area Under the ROC Curve (AUC)

- The area under a classifier's curve expressed as a fraction of the unit square
 - Its value ranges from zero to one
- The AUC is useful when a single number is needed to summarize performance, or when nothing is known about the operating conditions
 - A ROC curve provides more information than its area
- Equivalent to the Mann-Whitney-Wilcoxon measure
 - Also equivalent to the Gini Coefficient (with a minor algebraic transformation)
 - Both are equivalent to the probability that a randomly chosen positive instance will be ranked ahead of a randomly chosen negative instance

Performance evaluation

- Training Set:

Model	Accuracy
Classification Tree	95%
Logistic Regression	93%
k-Nearest Neighbors	100%
Naive Bayes	76%

- Test Set:

Model	Accuracy	AUC
Classification Tree	$91.8\% \pm 0.0$	0.614 ± 0.014
Logistic Regression	$93.0\% \pm 0.1$	0.574 ± 0.023
k-Nearest Neighbors	$93.0\% \pm 0.0$	0.537 ± 0.015
Naive Bayes	$76.5\% \pm 0.6$	0.632 ± 0.019

Performance evaluation

- Naive Bayes confusion matrix:

	p	n
1	127 (3%)	200 (4%)
0	848 (18%)	3518 (75%)

- k -Nearest Neighbors confusion matrix:

	p	n
1	3 (0%)	324 (7%)
0	15 (0%)	4351 (93%)

ROC Curve

