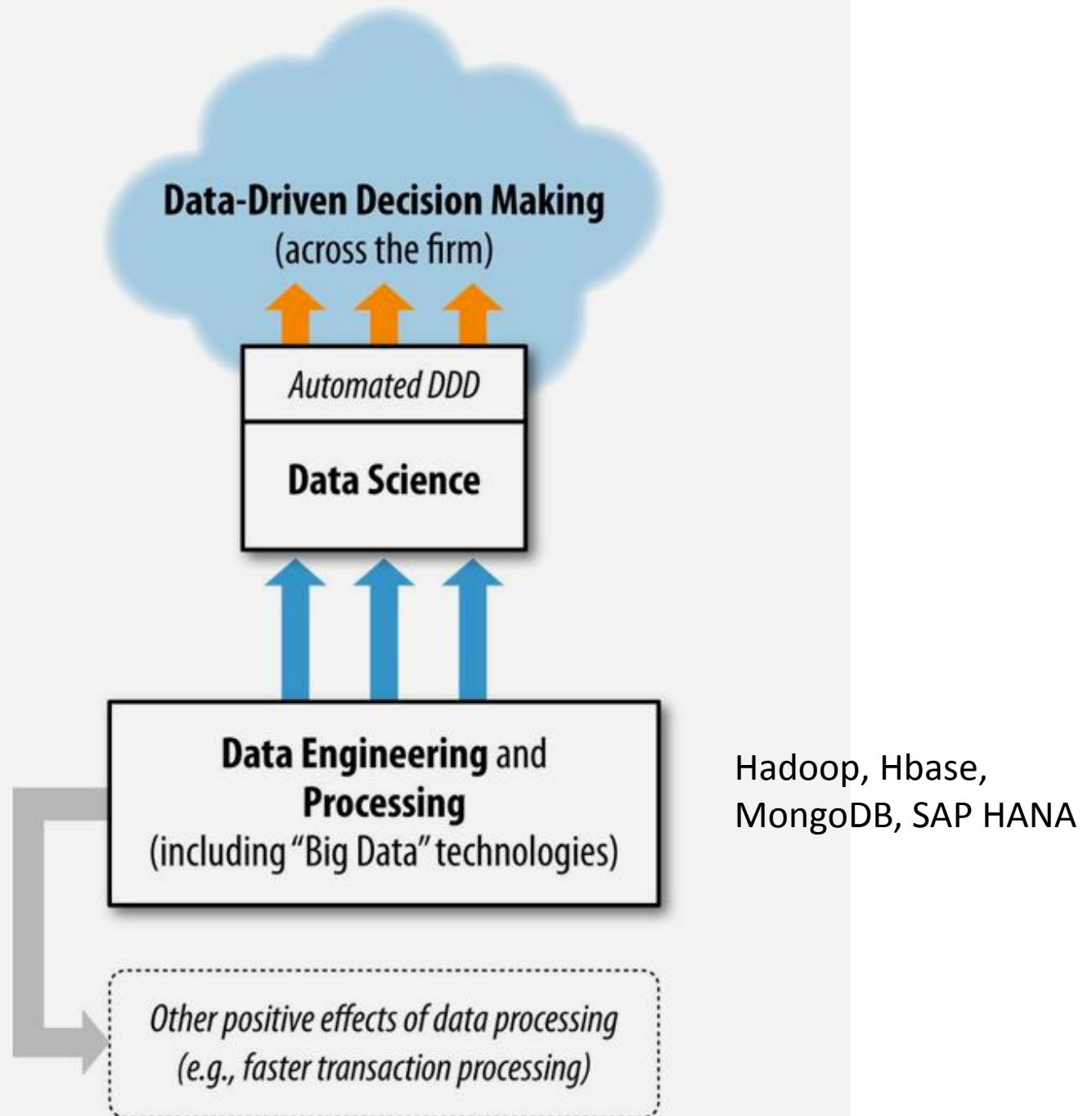


Data Science for Business

陳炫碩

中央大學 ERP中心 主任

Extracting knowledge from data to
solve business problems



Data Science & Data Mining

- Data science is a set of fundamental principles that guide the extraction of knowledge from data.
- Data mining is the extraction of knowledge from data, via technologies that incorporate these principles.

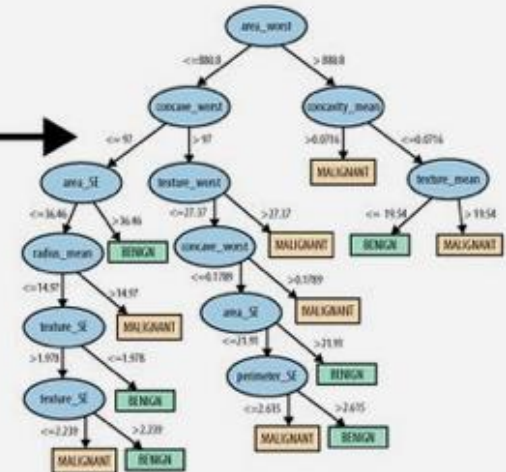
Historical Data

x	y	z	class
14	True	Red	accepted
6	True	Blue	rejected
...			
50.3	False	Red	accepted

Data mining



Model



Training data have all values specified

Model is deployed

Mining

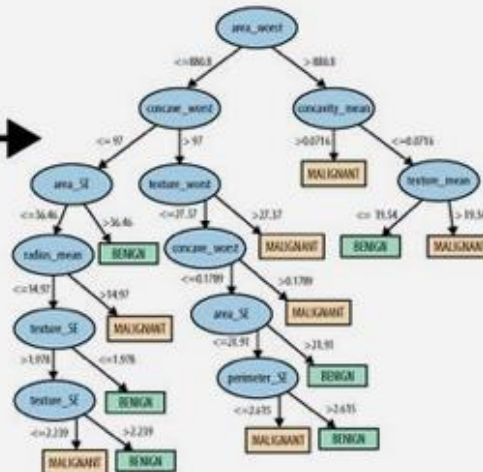
Use

New data item

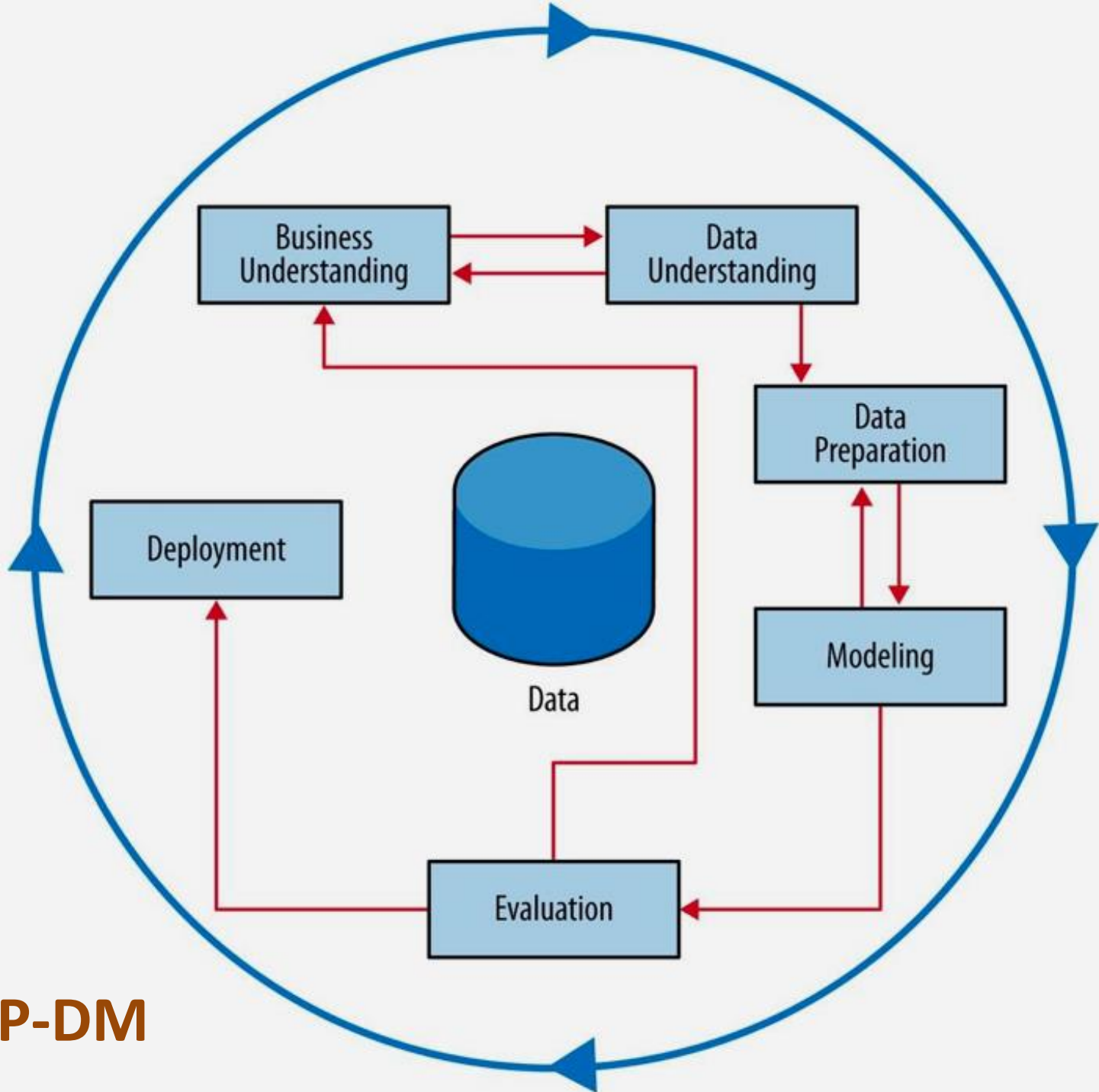
x	y	z	class
30	false	Red	?

New data item has class value unknown (e.g. will customer accept?)

Model



**Class: accepted,
Probability: 0.88**



CRISP-DM

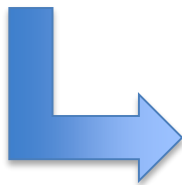
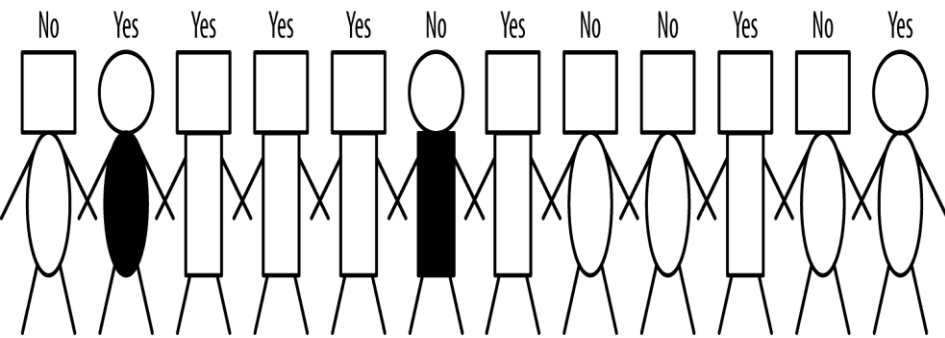
Modeling

Predictive model

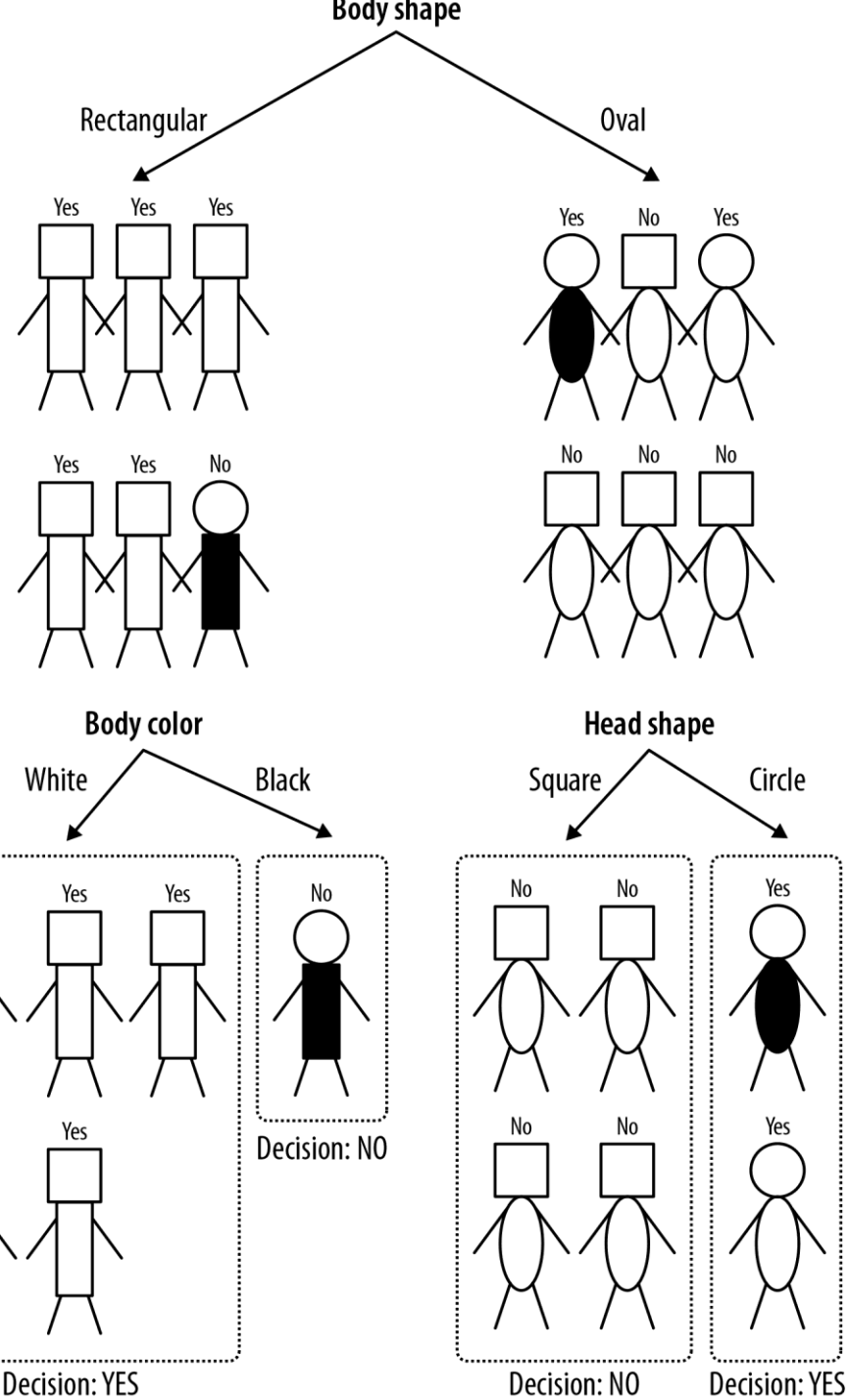
- Informative variables
 - provide information that reduces uncertainty about something
- Model: simplified representation of reality created to serve a purpose
- A predictive model is a formula for estimating the unknown value of interest. (the formula could be mathematical or a logical statement.)

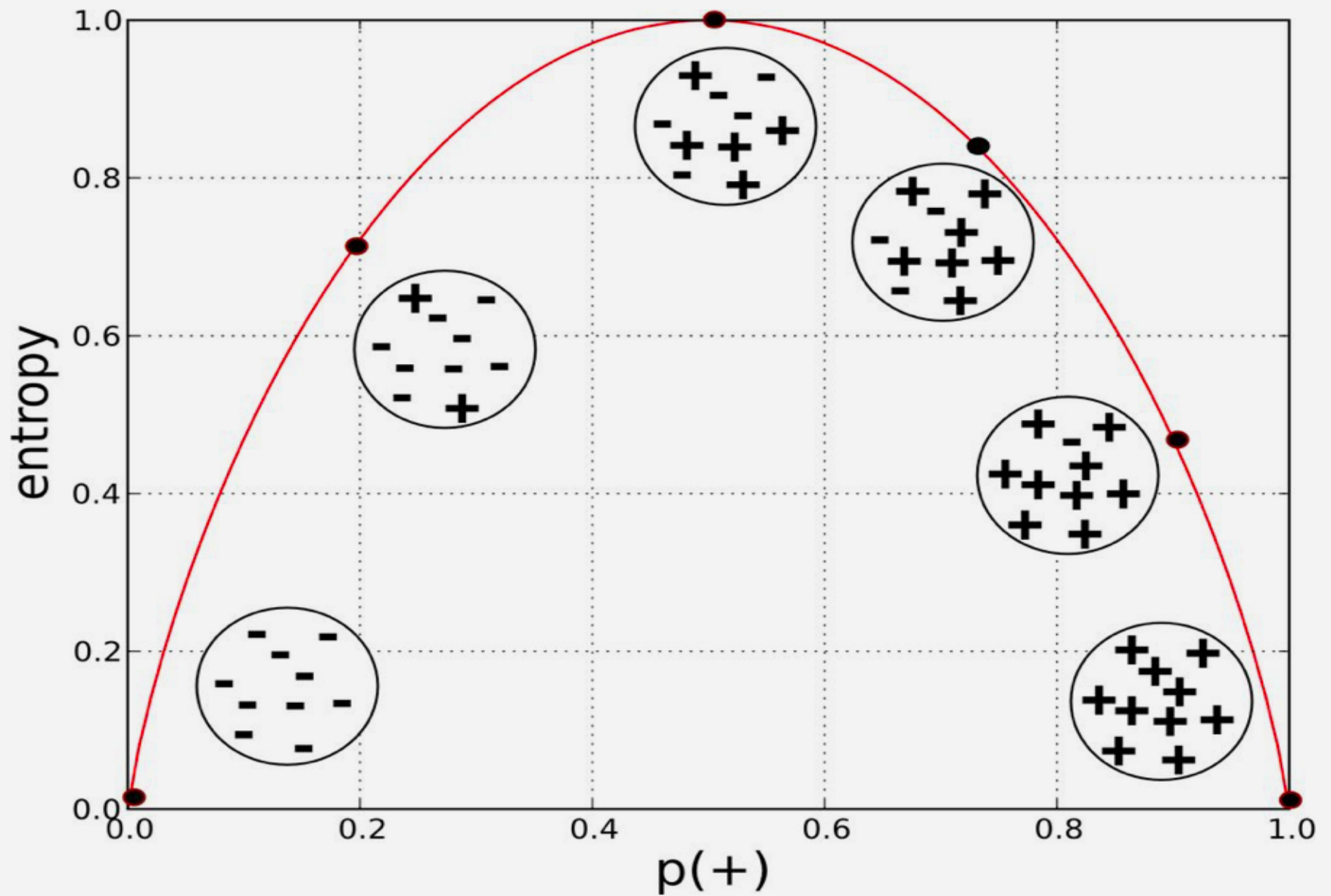
Prediction model

- Classification model
- Class-probability estimation model
- Regression model



Classification





$$\text{Entropy} = \sum_i -p_i \log_2 p_i$$

Information gain

- How much pure are the split sets compared to the original set?
- How much an attribute improves entropy?

Decision Tree

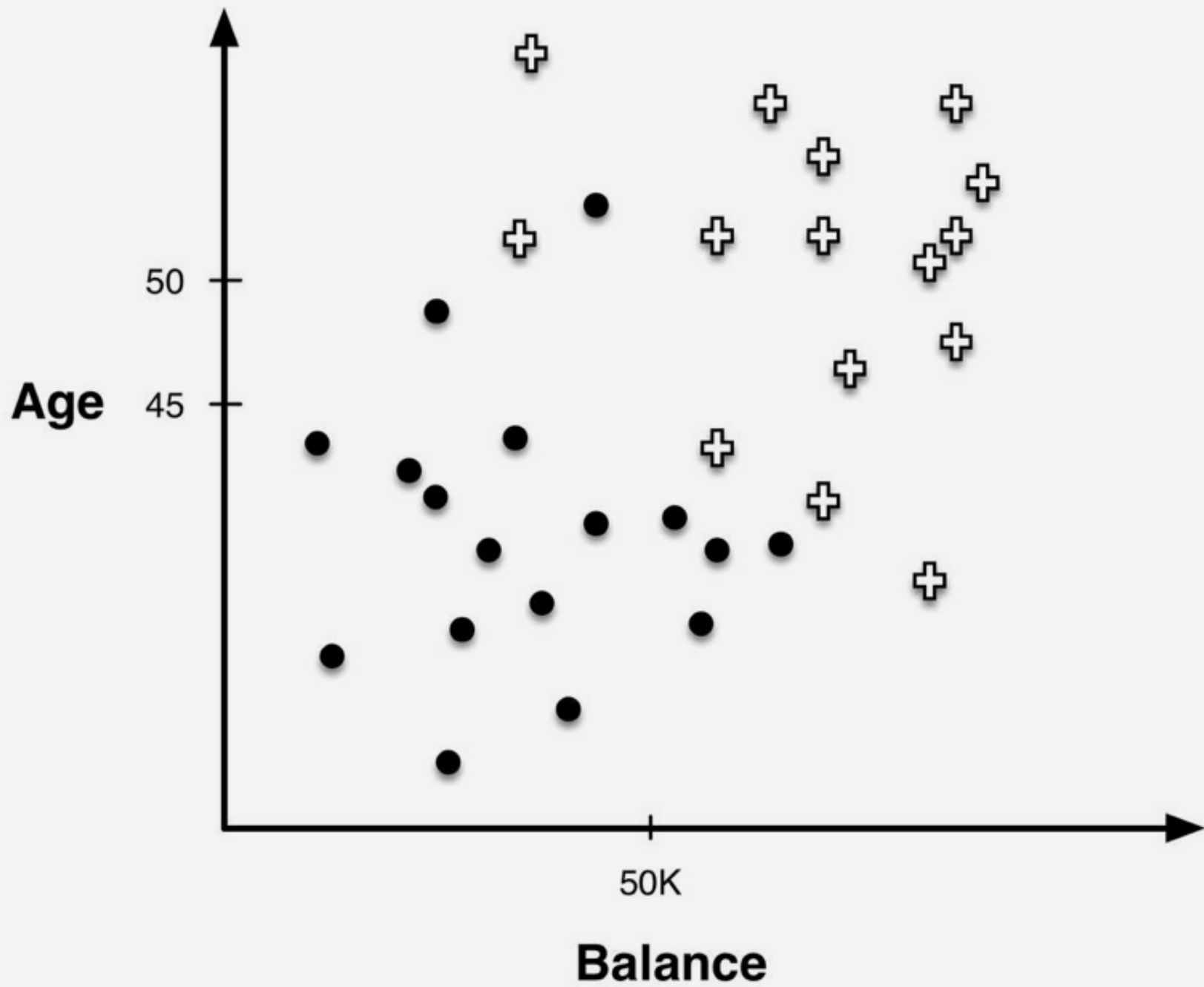
The diagram shows a table with 5 columns and 6 rows. The first four columns are grouped under the label 'Attributes' with a bracket above them. The fifth column is labeled 'Target attribute' with an arrow pointing to it. The third row (Claudio) is highlighted in blue. An arrow points from this row to the explanatory text below.

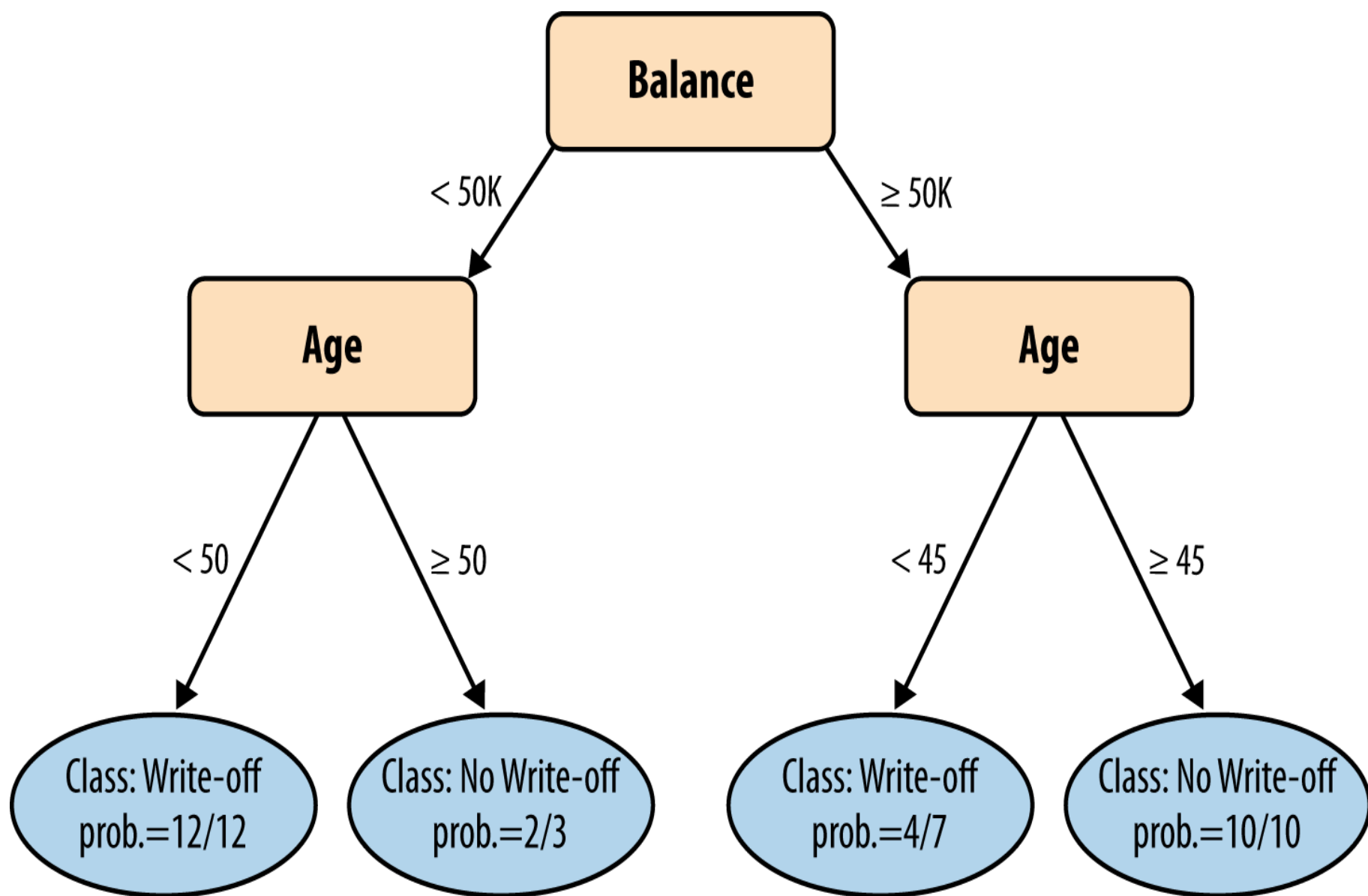
Attributes				Target attribute
Name	Balance	Age	Employed	Write-off
Mike	\$200,000	42	no	yes
Mary	\$35,000	33	yes	no
Claudio	\$115,000	40	no	no
Robert	\$29,000	23	yes	yes
Dora	\$72,000	31	no	no

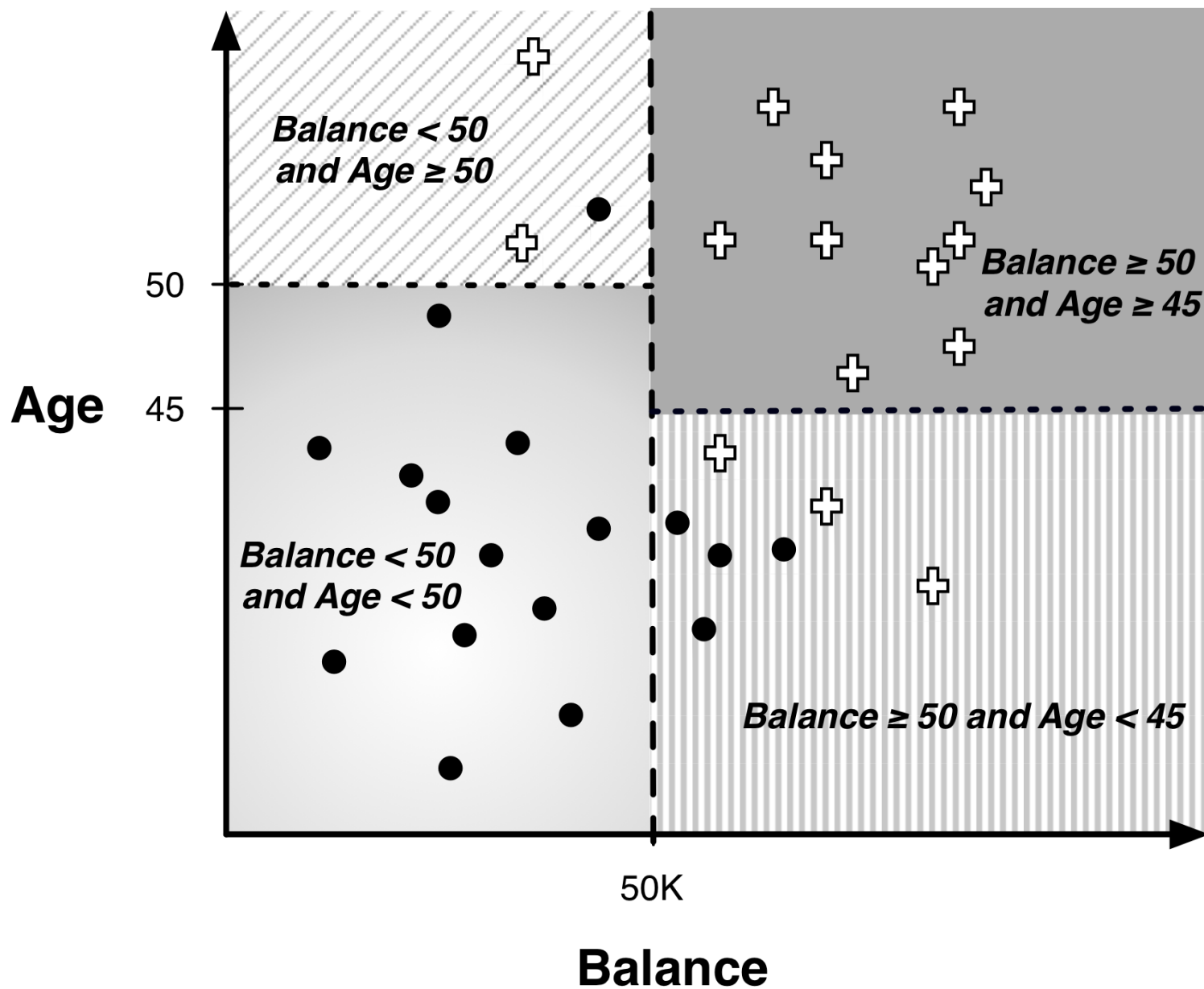
This is one row (example).

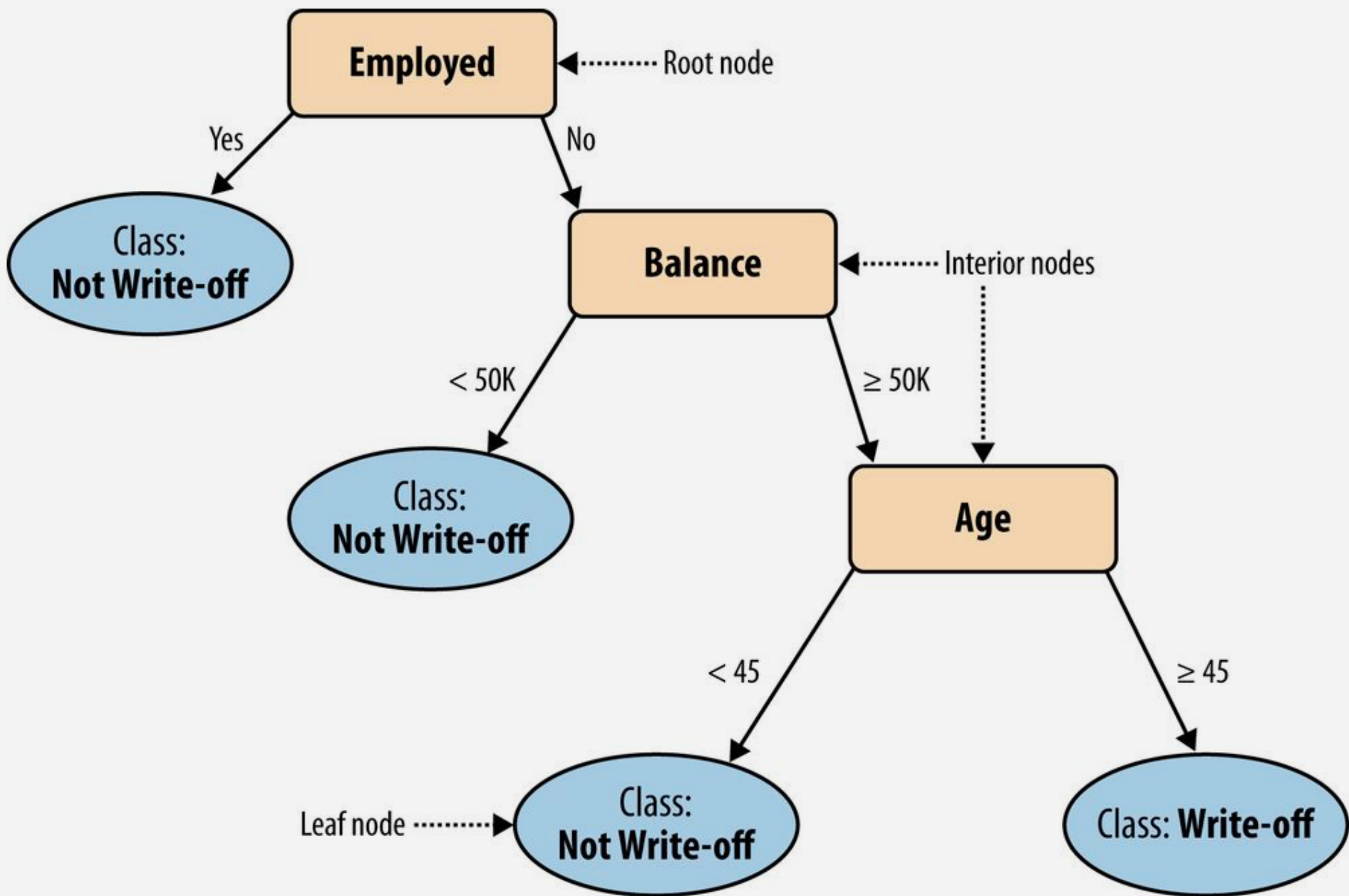
Feature vector is: **<Claudio,115000,40,no>**

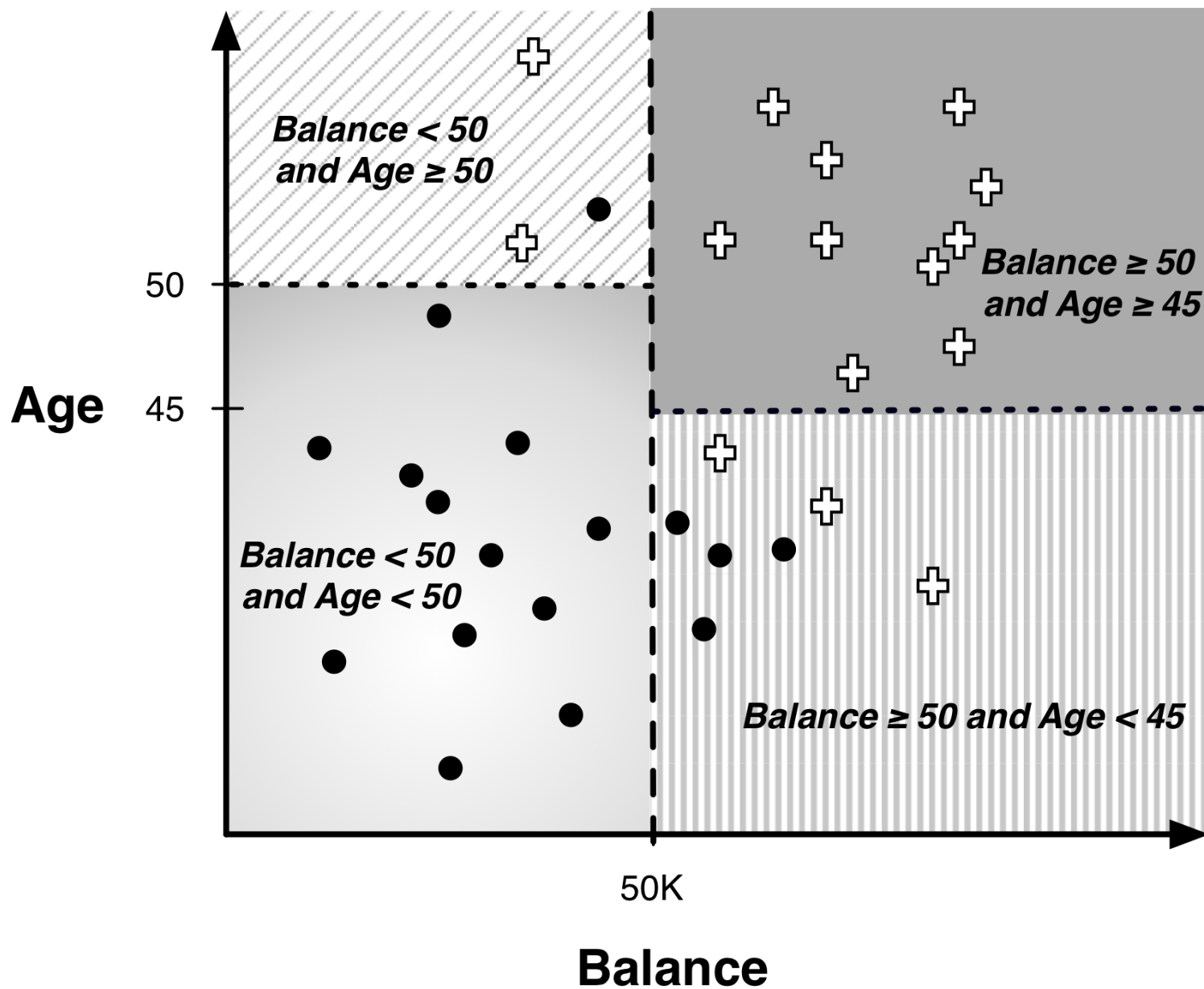
Class label (value of Target attribute) is **no**







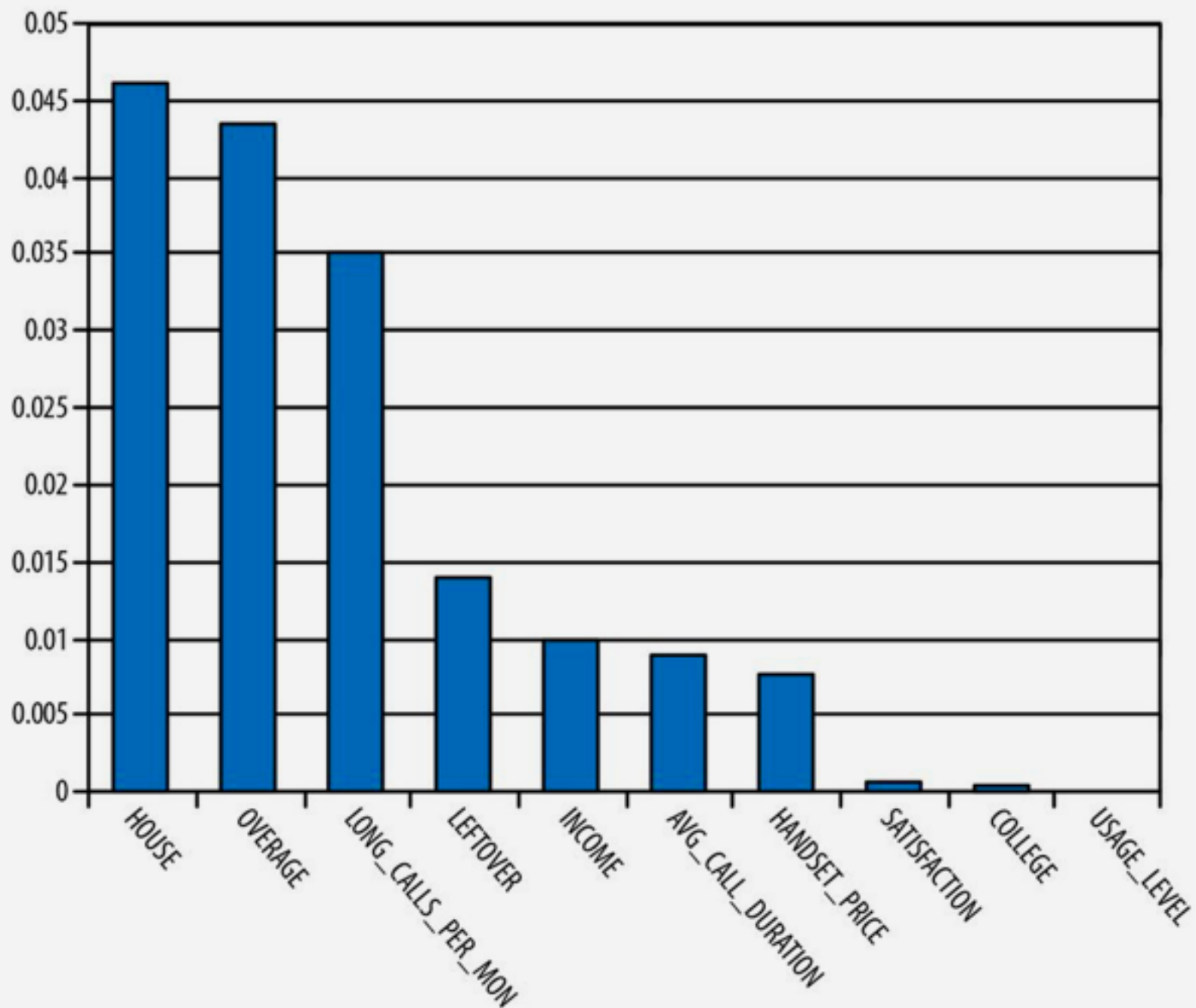


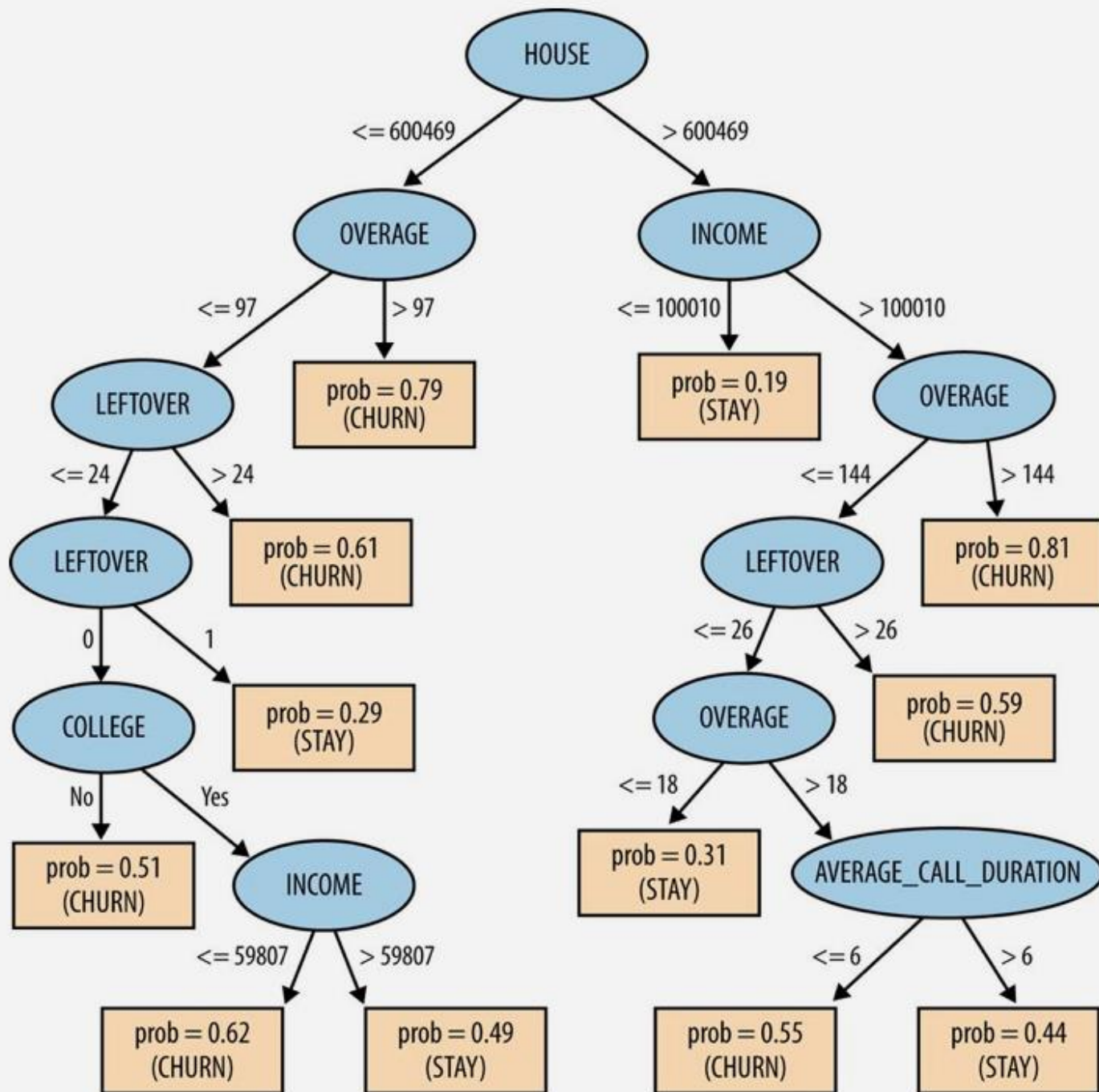


Class-probability estimation model

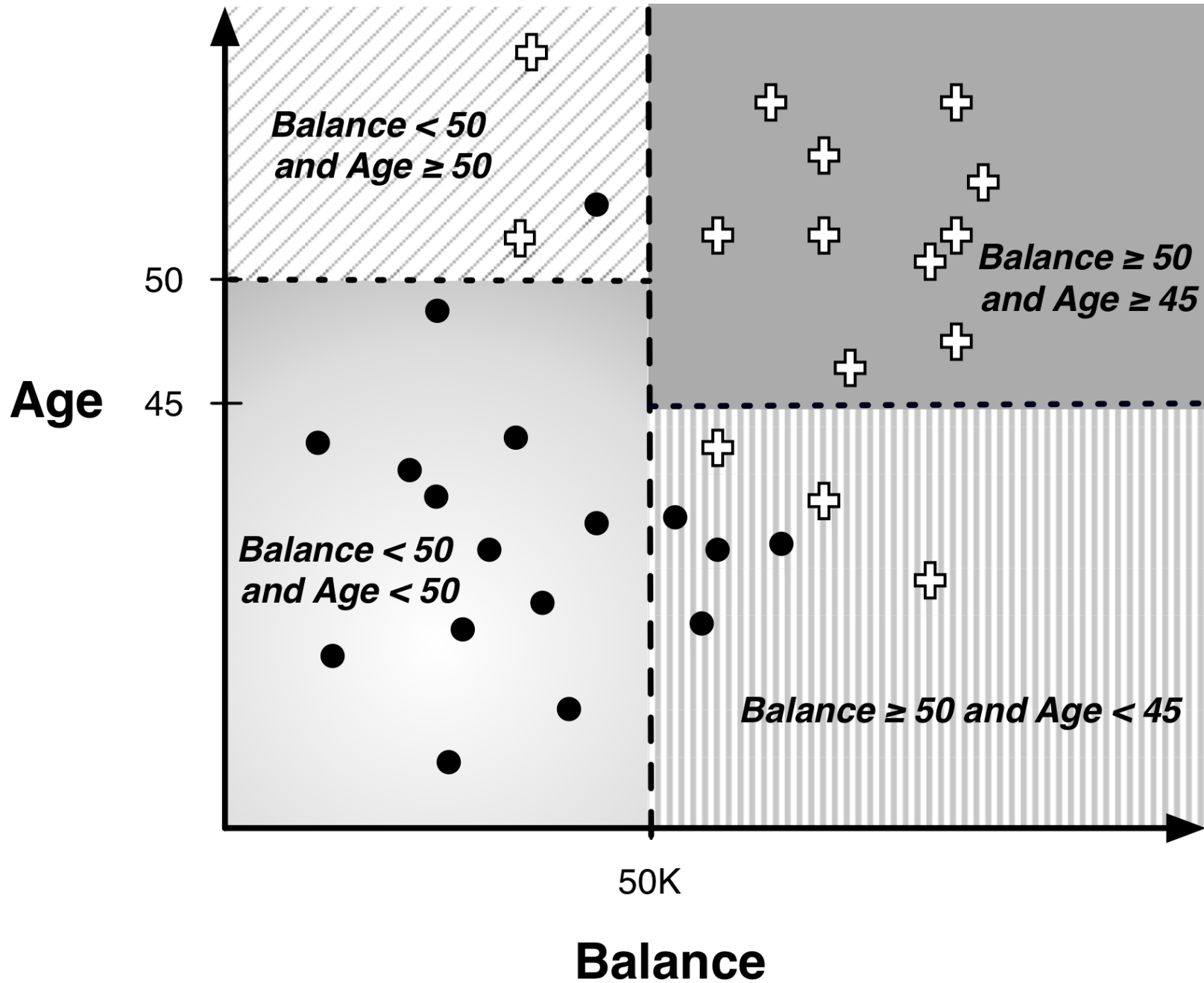
- Frequency-based estimate of class membership probability.
- Laplace correction: to moderate the influence of leaves with only a few instances.

$$p(c) = (n+1)/(n+m+2)$$



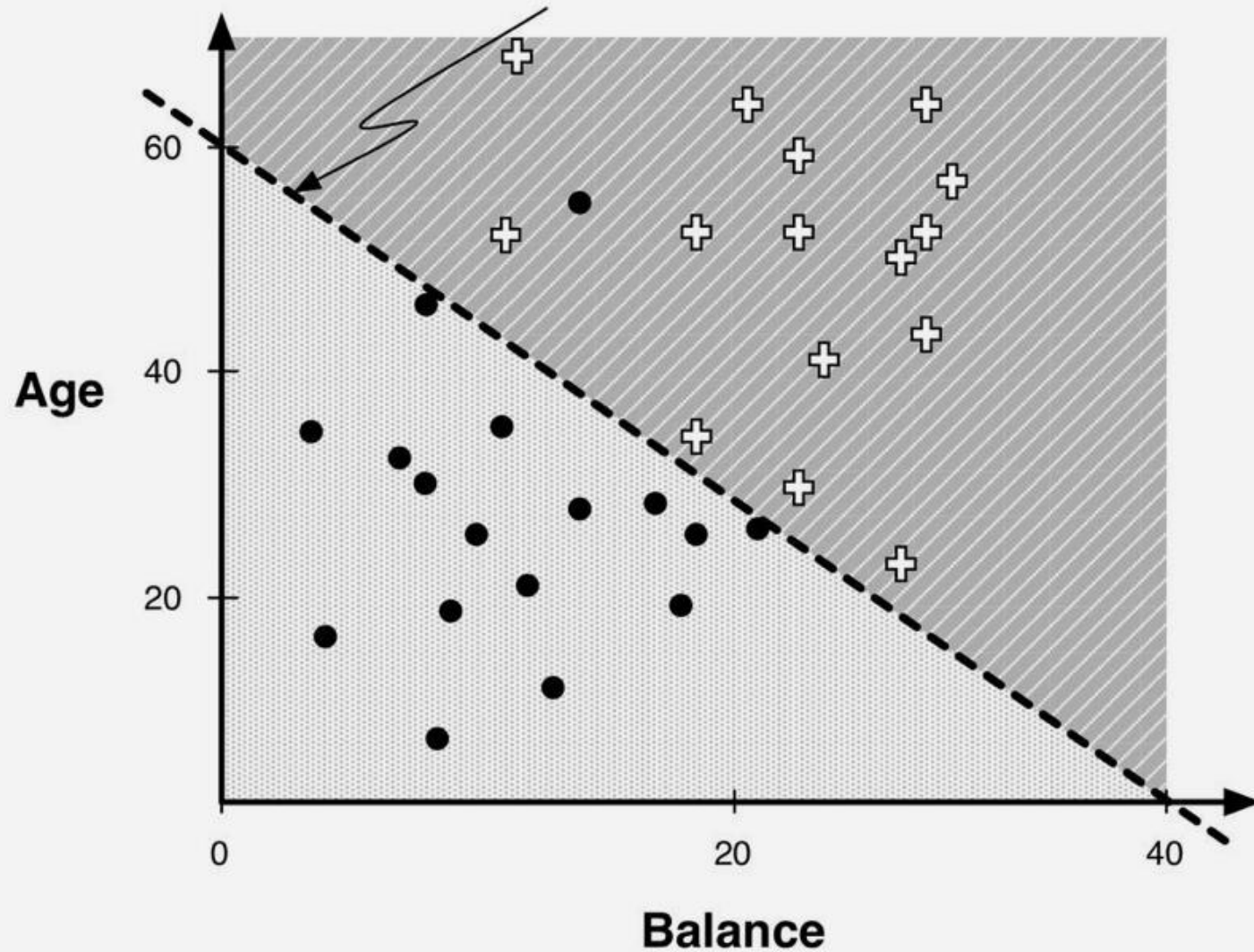


Linear Model



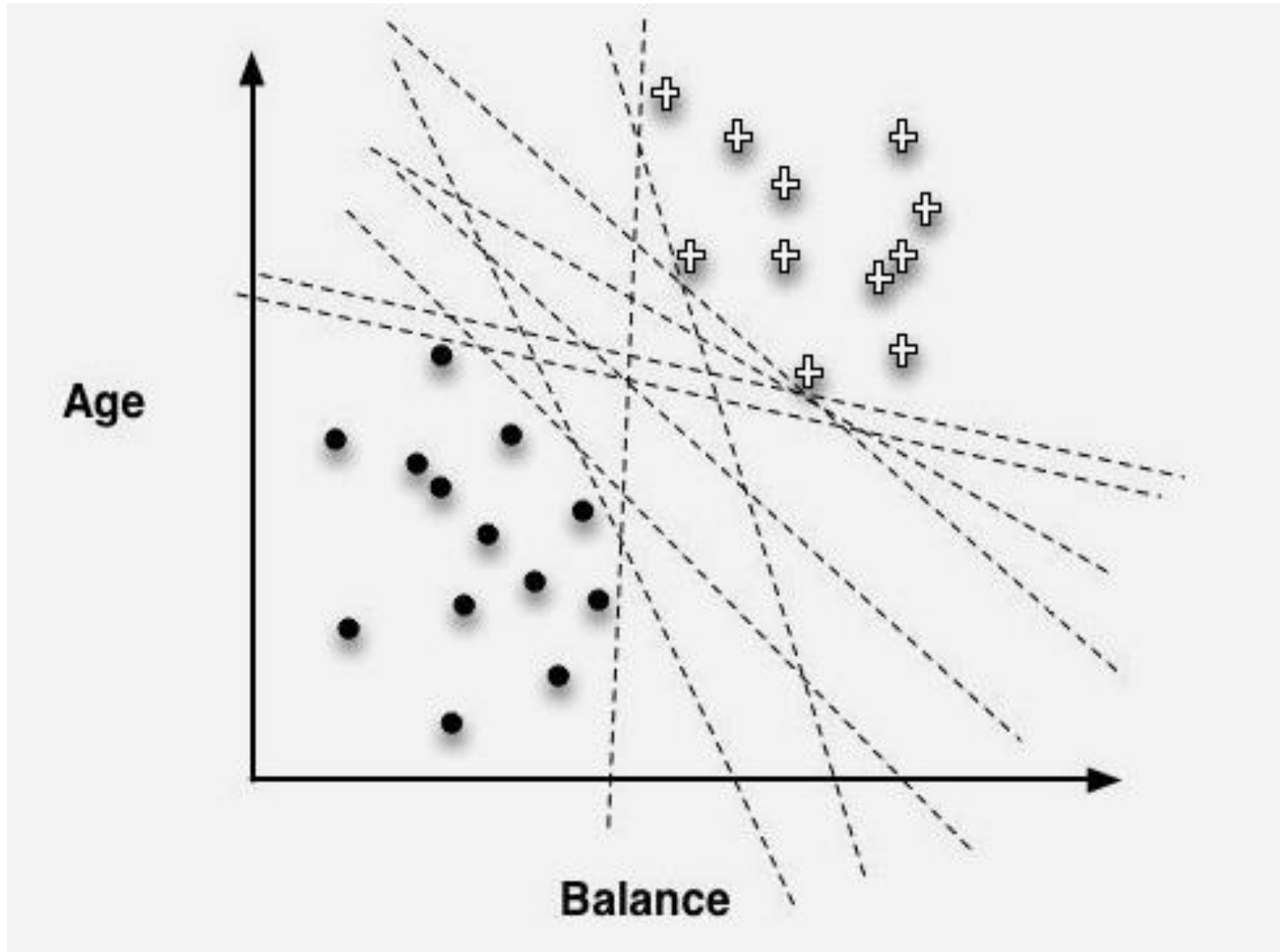
Decision boundary:

$$\text{Age} = \text{Balance} \times -1.5 + 60$$



$$f(\mathbf{x}) = 60 - 1.0 \times \text{Age} - 1.5 \times \text{Balance}$$

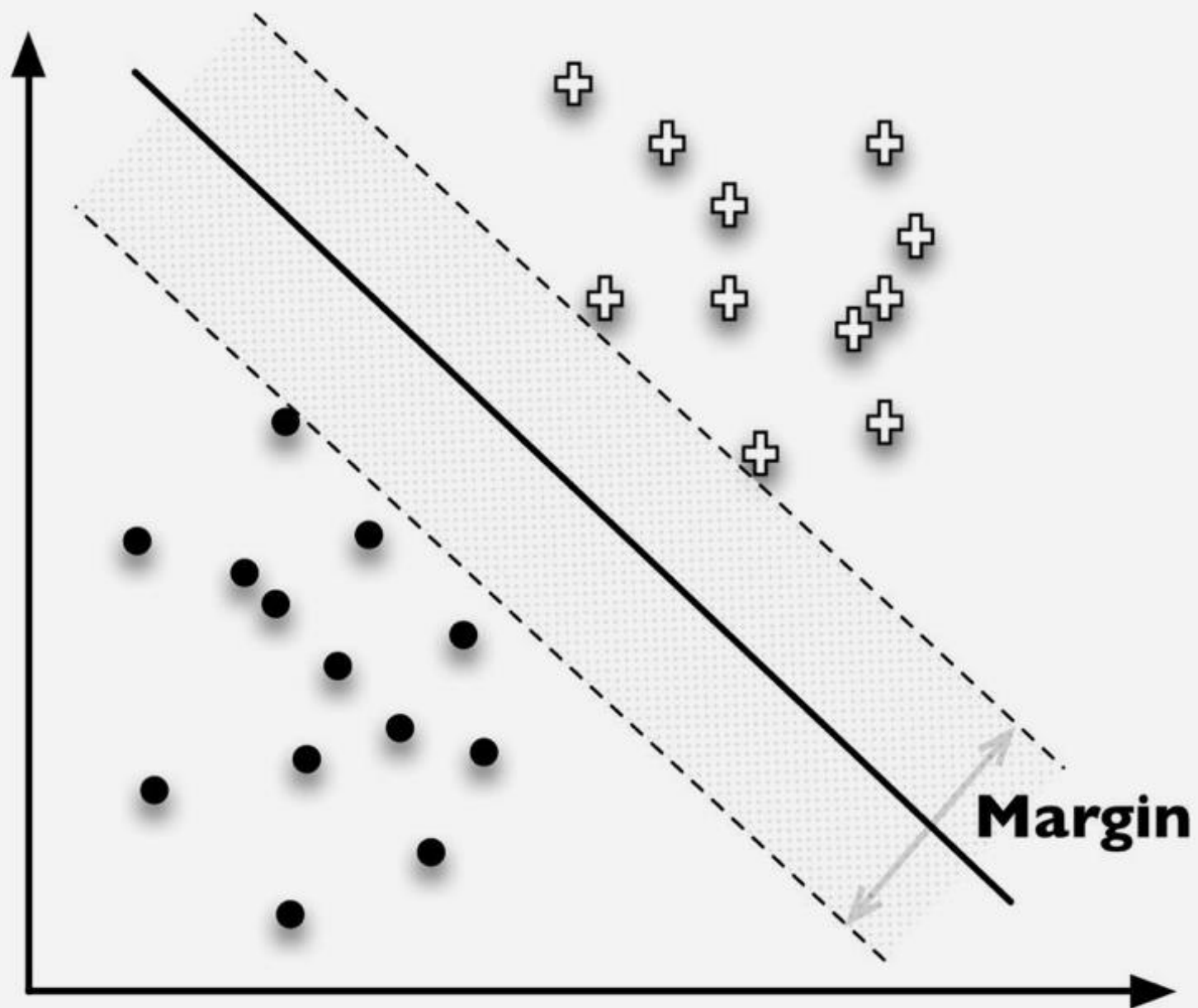
Which is the best line to fit the data?

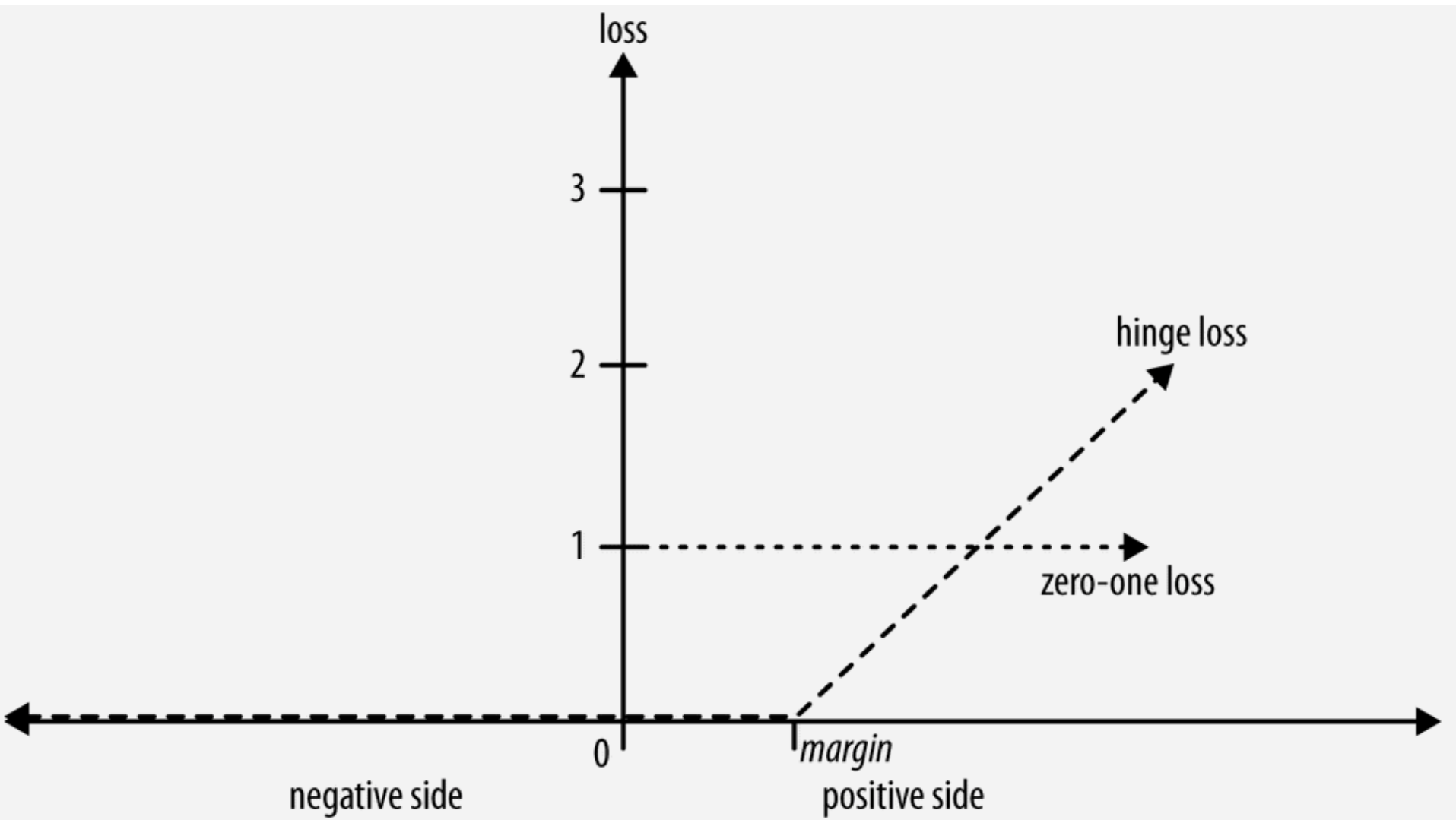


Objective function

- Linear regression, support vector machines (SVM) and logistic regression are all trying to fit a linear model to data.
- The key difference is that each uses a different objective function.
- Linear regression -> least square
- SVM -> Maximizing the margin
- Logistic regression -> Maximum likelihood

Age





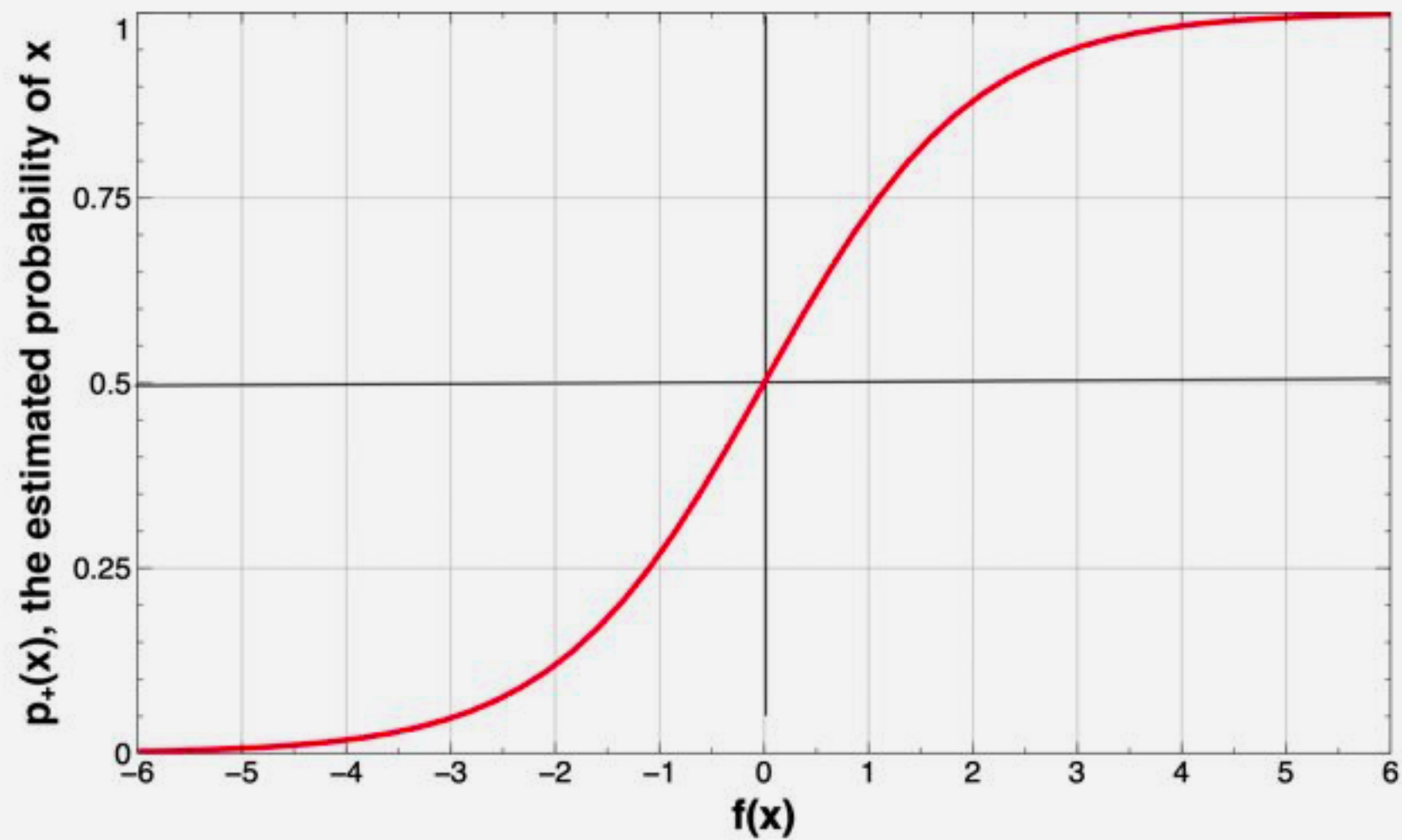
Class probability estimation

- What is the problem with simply using basic linear model to estimate the class probability?
- $f(x)$ gives the distance from the separating boundary, ranging from $-\infty$ to $+\infty$, while a probability should range from 0 to 1.

Probability	Odds	Log-odds
0.5	50:50 or 1	0
0.9	90:10 or 9	2.19
0.999	999:1 or 999	6.9
0.01	1:99 or 0.0101	-4.6
0.001	1:999 or 0.001001	-6.9

$$\log \left(\frac{p_+(\mathbf{x})}{1 - p_+(\mathbf{x})} \right) = f(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + \dots$$

$$p_+(\mathbf{x}) = \frac{1}{1 + e^{-f(\mathbf{x})}}$$



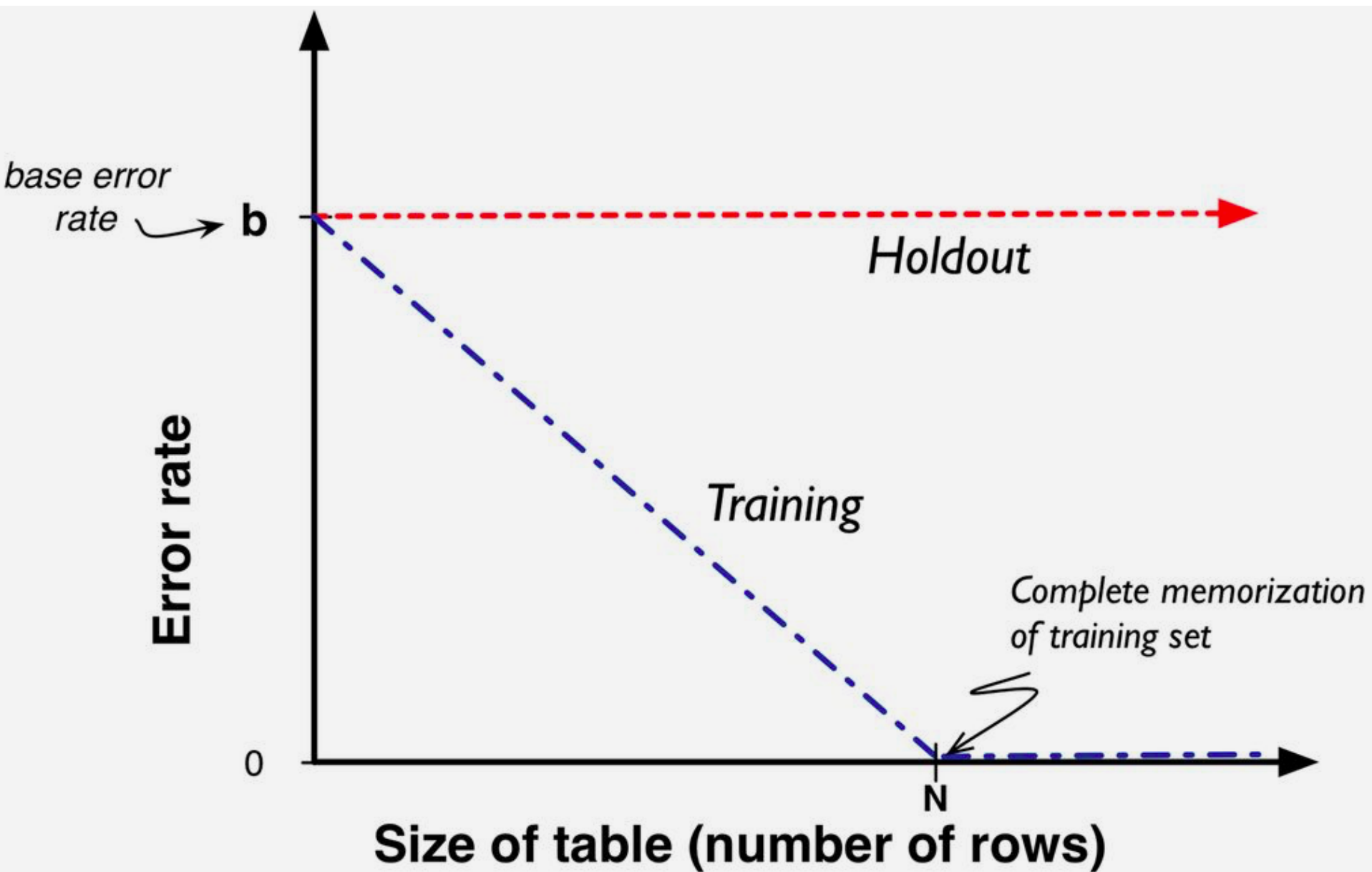
Objective function for Logistic Regression

$$g(\mathbf{x}, \mathbf{w}) = \begin{cases} p_+(\mathbf{x}) & \text{if } \mathbf{x} \text{ is a } + \\ 1 - p_+(\mathbf{x}) & \text{if } \mathbf{x} \text{ is a } \bullet \end{cases}$$

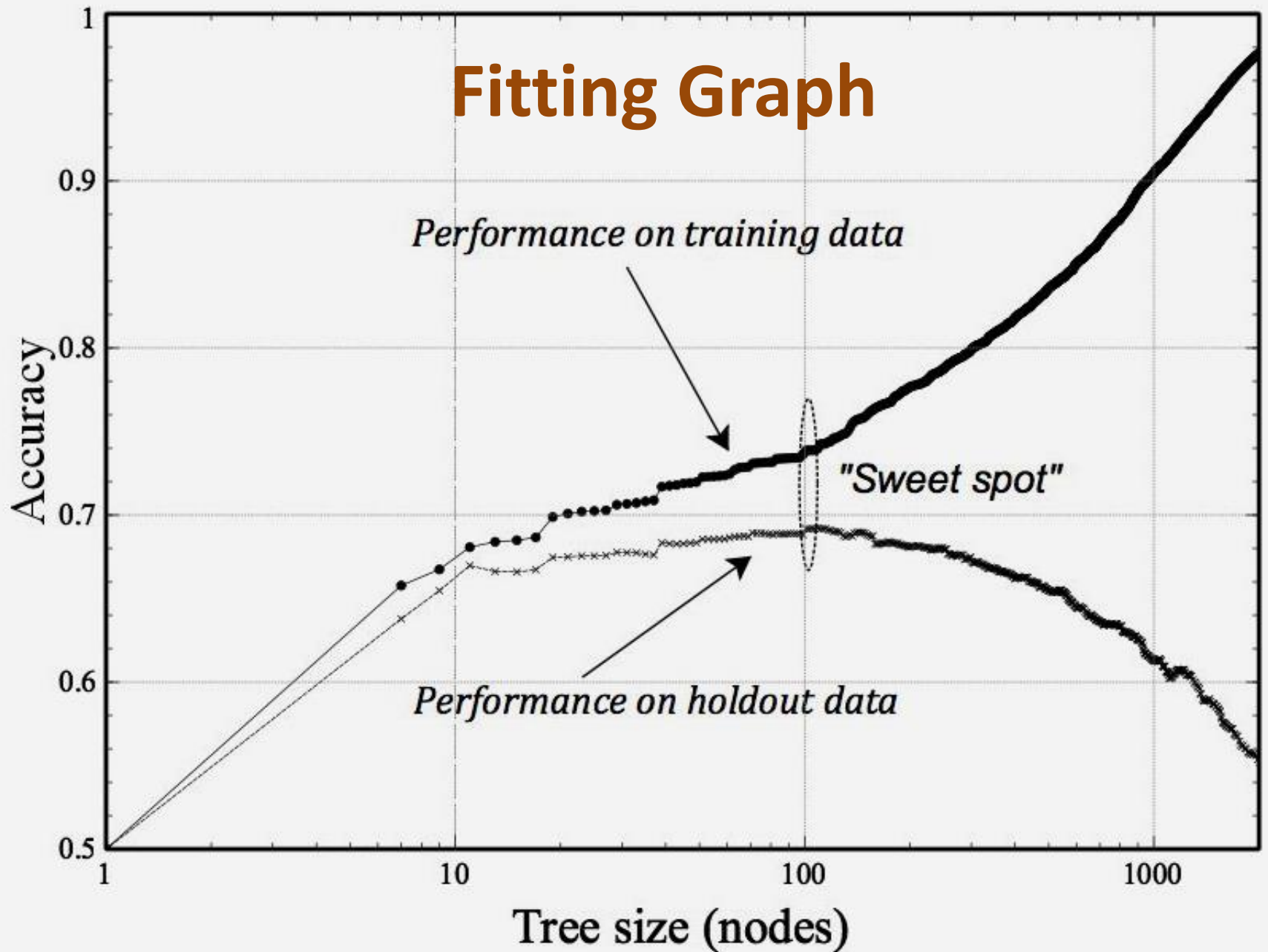
Overfitting

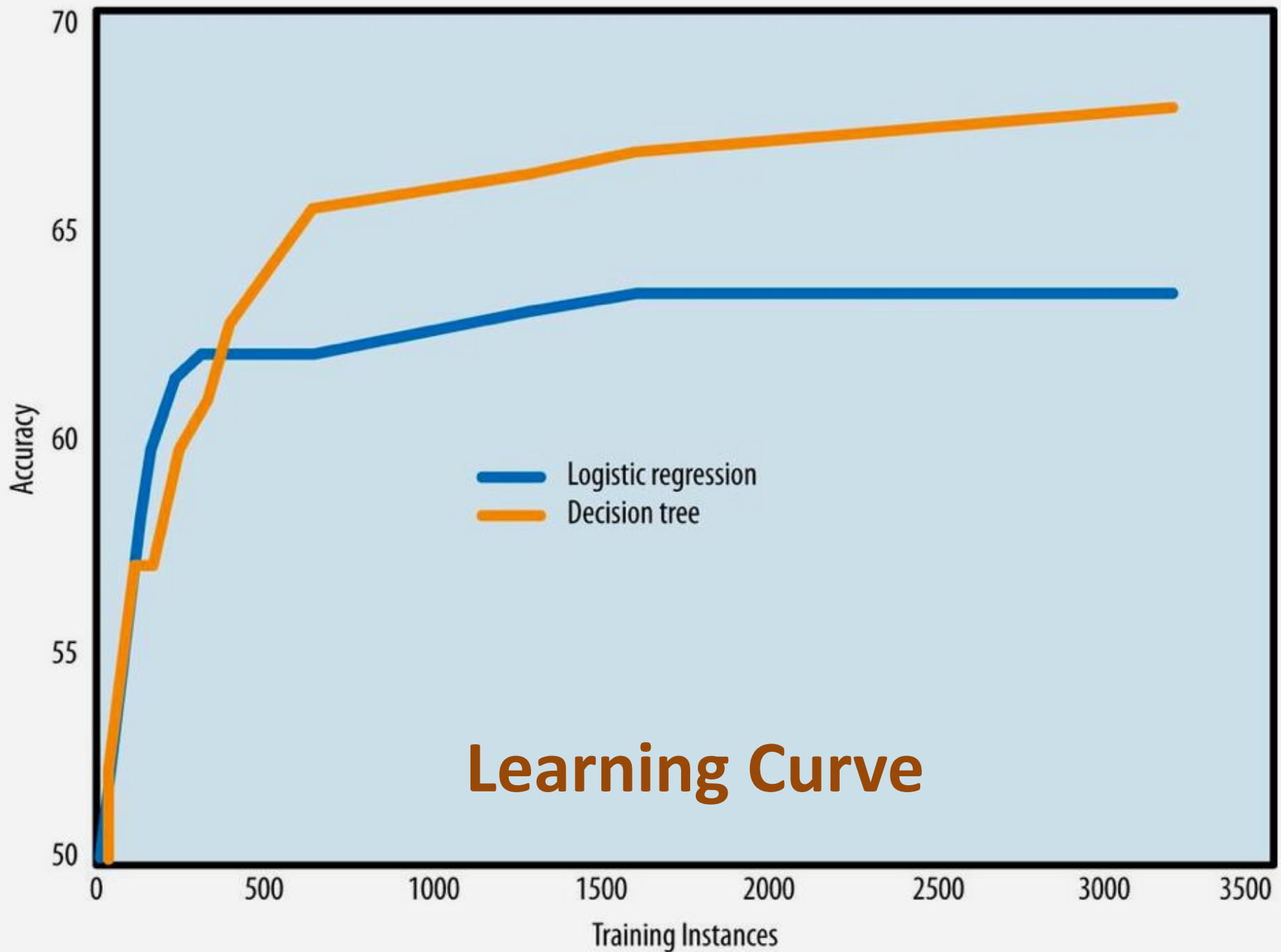
Overfitting

- If we allow ourselves enough flexibility in searching for patterns in a particular dataset, we will find patterns.
- These “patterns” may be just chance occurrences in the data, while we are interested in patterns that generalize – that predict well for instances that we have not yet observed.
- For example, a table model memorizes the training data and performs no generalization.

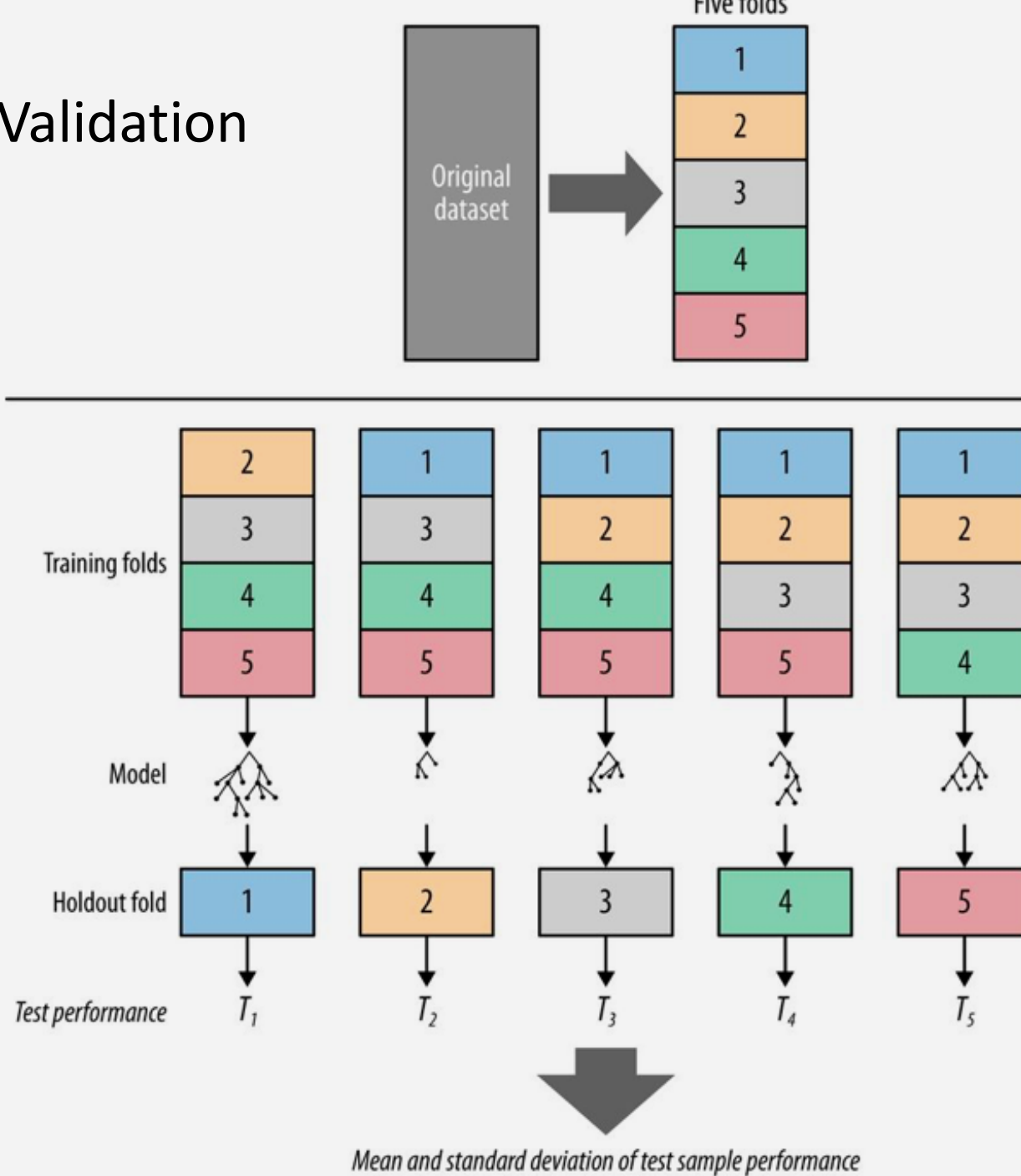


Fitting Graph





Cross-Validation



What is a good model?

Evaluation

- Classifiers
- Ordering
- Class probability

Bad Positive & Harmless Negative

- Positive: worthy of attention
- Negative: uninteresting
- A classifier is screening through a large population consisting mostly of negatives and looking for a small number of positive instances.
- The number of mistakes made on negative examples (false positive) may be dominant, though the cost of each mistakes made on positive example (false negative) will be higher.

Accuracy?

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

Confusion Matrix

	p	n
Y	True positives	False positives
N	False negatives	True negatives

Problems with unbalanced classes

- A Tele company has churn rate of 0.01%.
- Is a model with 95% accuracy good?

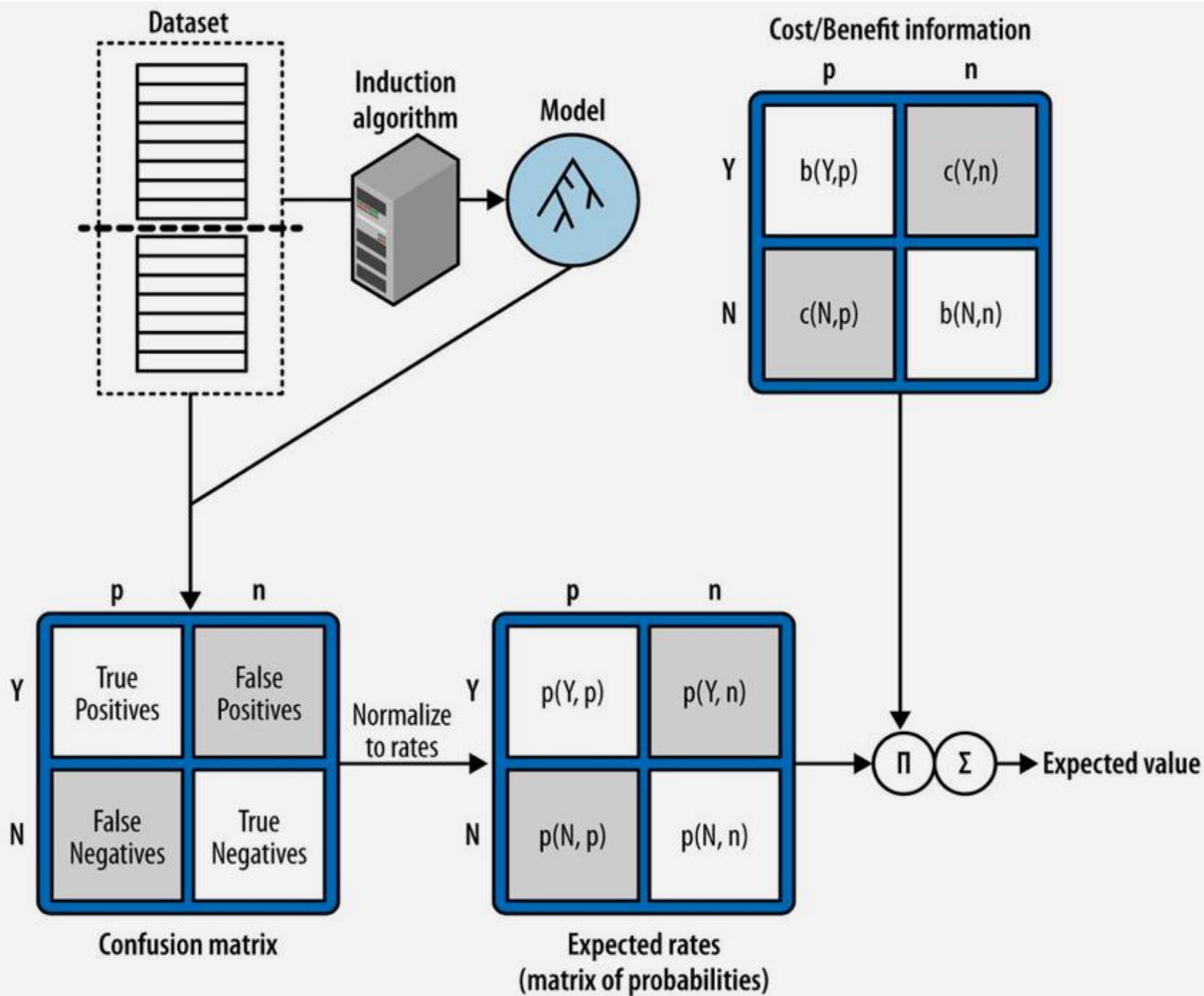
Problems with unequal costs and benefits

- How much we care about different errors? (false negatives & false positives)
- What are the good results for the correct prediction? (true positives & true negatives)
- The expected value is the weighted average of the values from different possible outcomes where the weight to each value is the probability of occurrence.

Example

- Promotion cost: \$1
- Product-related cost: \$100
- Product price: \$150

Who are the customers this campaign should target to make a profit?



	p	n
Y	56	7
N	5	42

		Actual	
		p	n
Predicted	Y	b(Y,p)	c(Y,n)
	N	c(N,p)	b(N,n)

$$T = 110$$

$$p(\mathbf{Y}, \mathbf{p}) = 56/110 = 0.51$$

$$p(\mathbf{Y}, \mathbf{n}) = 7/110 = 0.06$$

$$p(\mathbf{N}, \mathbf{p}) = 5/110 = 0.05$$

$$p(\mathbf{N}, \mathbf{n}) = 42/110 = 0.38$$

$$\text{Expected profit} = p(\mathbf{Y}, \mathbf{p}) \cdot b(\mathbf{Y}, \mathbf{p}) + p(\mathbf{N}, \mathbf{p}) \cdot b(\mathbf{N}, \mathbf{p}) + \\ p(\mathbf{N}, \mathbf{n}) \cdot b(\mathbf{N}, \mathbf{n}) + p(\mathbf{Y}, \mathbf{n}) \cdot b(\mathbf{Y}, \mathbf{n})$$

$$p(x, y) = p(y) \cdot p(x \mid y)$$

$$\text{Expected profit} = p(\mathbf{Y} \mid \mathbf{p}) \cdot p(\mathbf{p}) \cdot b(\mathbf{Y}, \mathbf{p}) + p(\mathbf{N} \mid \mathbf{p}) \cdot p(\mathbf{p}) \cdot b(\mathbf{N}, \mathbf{p}) + \\ p(\mathbf{N} \mid \mathbf{n}) \cdot p(\mathbf{n}) \cdot b(\mathbf{N}, \mathbf{n}) + p(\mathbf{Y} \mid \mathbf{n}) \cdot p(\mathbf{n}) \cdot b(\mathbf{Y}, \mathbf{n})$$

$$\text{Expected profit} = p(\mathbf{p}) \cdot [p(\mathbf{Y} \mid \mathbf{p}) \cdot b(\mathbf{Y}, \mathbf{p}) + p(\mathbf{N} \mid \mathbf{p}) \cdot b(\mathbf{N}, \mathbf{p})] + \\ p(\mathbf{n}) \cdot [p(\mathbf{N} \mid \mathbf{n}) \cdot b(\mathbf{N}, \mathbf{n}) + p(\mathbf{Y} \mid \mathbf{n}) \cdot b(\mathbf{Y}, \mathbf{n})]$$

Correct & Error Rate

- True positive rate
- True negative rate
- False positive rate
- False negative rate

Correct & Error Rate

When the instance is actually positive

- True positive rate
- False negative rate

When the instance is actually negative

- True negative rate
- False positive rate

$T = 110$	
$P = 61$	$N = 49$
$p(\mathbf{p}) = 0.55$	$p(\mathbf{n}) = 0.45$
$tp\ rate = 56/61 = 0.92$	$fp\ rate = 7/49 = 0.14$
$fn\ rate = 5/61 = 0.08$	$tn\ rate = 42/49 = 0.86$

Sensitivity & Specificity

- Sensitivity: 有病者被檢出為有病的機率 (true positive rate)

有病者: $TP + FN$

Sensitivity 高 = FN 低, 即 false negative 少, 只要是檢驗結果是negative, 沒病機率高

- Specificity: 沒病者被檢出為沒病的機率

沒病者: $TN + FP$

Specificity 高 = FP低, 即 false positive 少, 只要是檢驗結果是positive, 有病機率高

Evaluation

- Classifiers
- Ordering
- Class probability

Instance description	True class	Score
.....	p	0.99
.....	p	0.98
.....	n	0.96
.....	n	0.90
.....	p	0.88
.....	n	0.87
.....	p	0.85
.....	p	0.80
.....	n	0.70
.....	p	0.65
.....	.	.
.....	.	.
.....	.	.

Y

p	n
0	0
N	100
	100

Y

p	n
1	0
N	99
	100

Y

p	n
2	0
N	98
	100

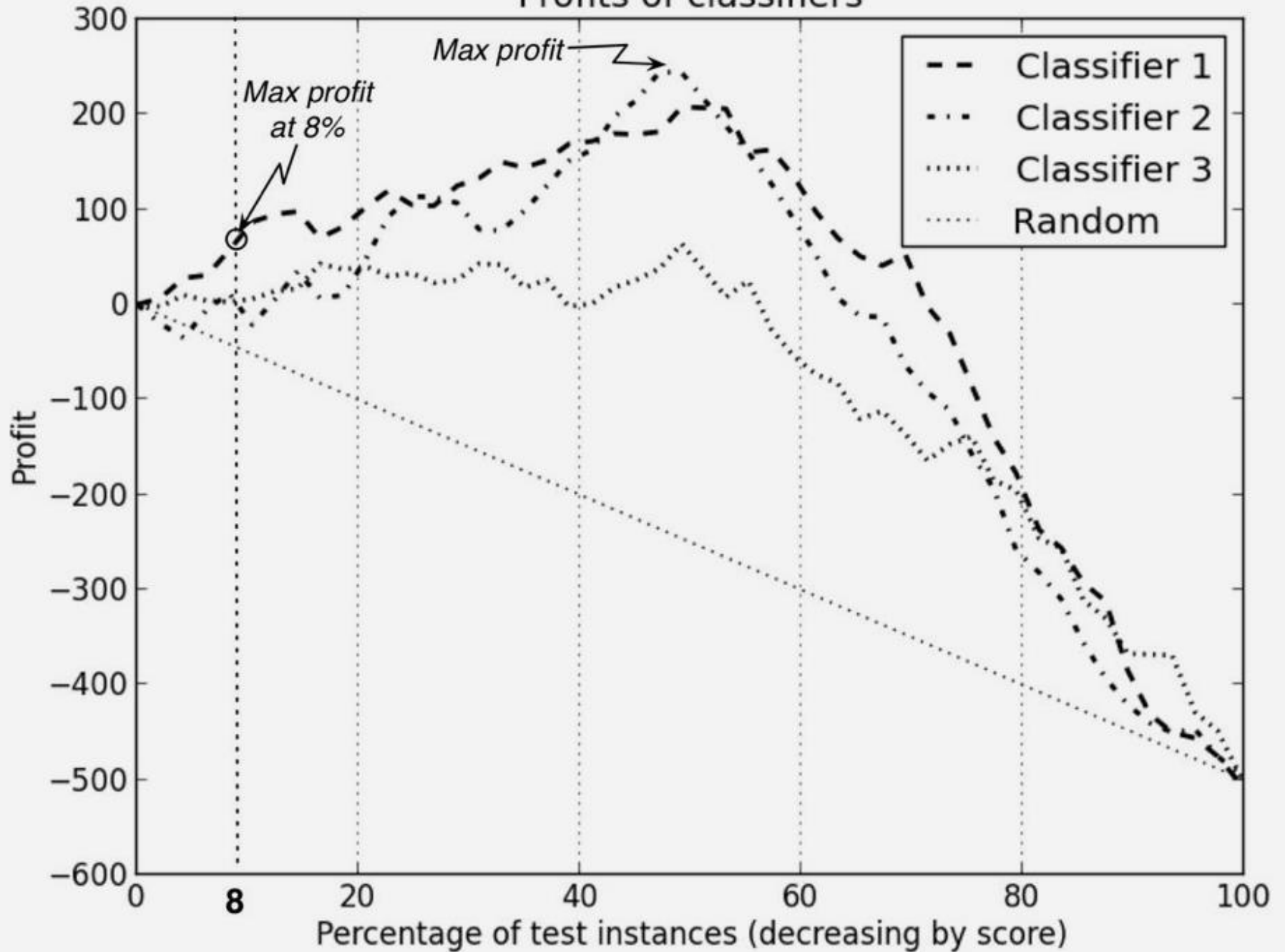
Y

p	n
2	1
N	98
	99

Y

p	n
6	4
N	94
	96

Profits of classifiers



Which model do you choose?

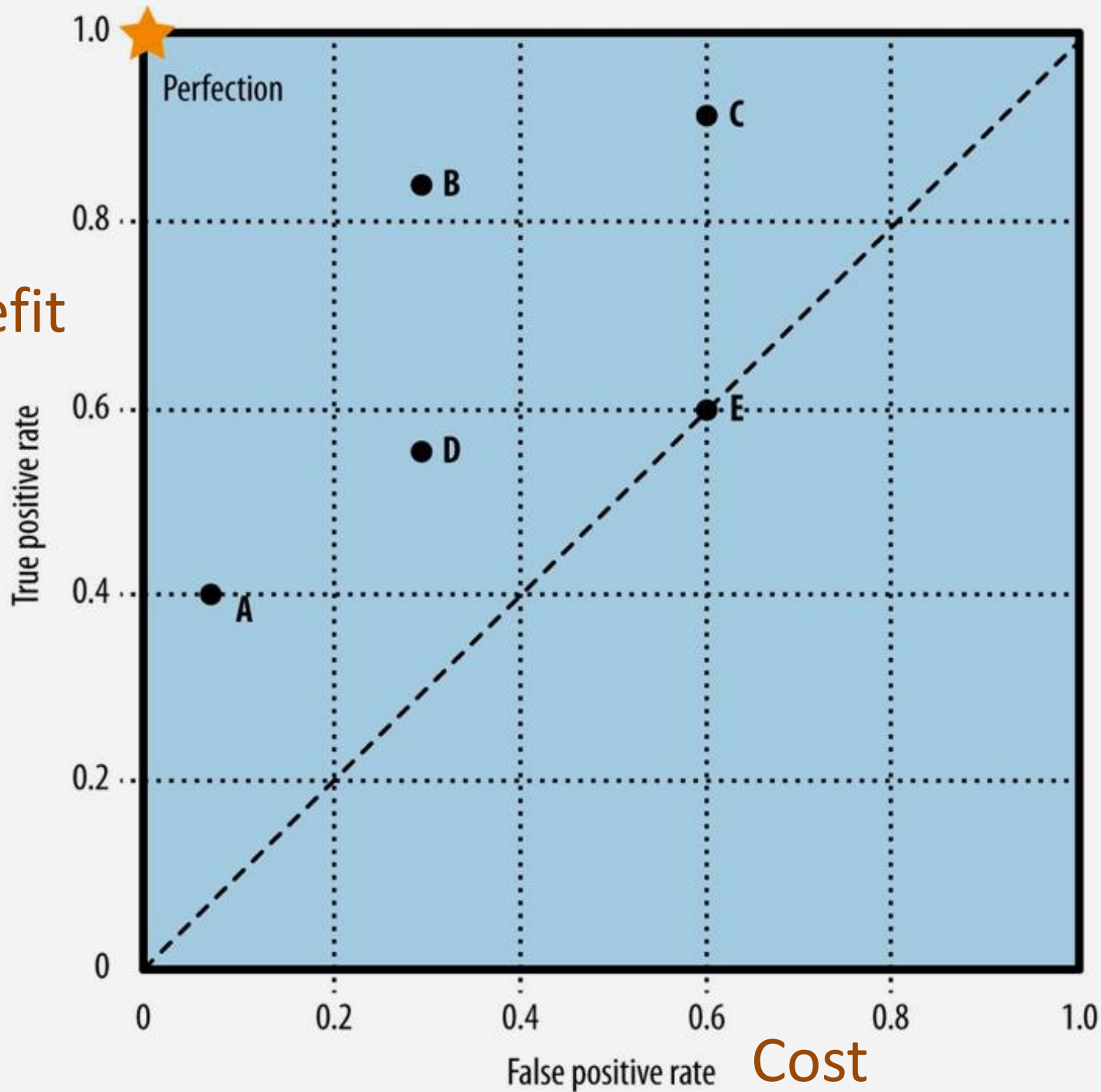
- If you have unlimited budget
- If you have 100,000 customers and \$40,000 budget (each offer costs \$5)

Limits of profit curve

- The class priors $p(p)$ & $p(n)$
- The costs and benefits

ROC

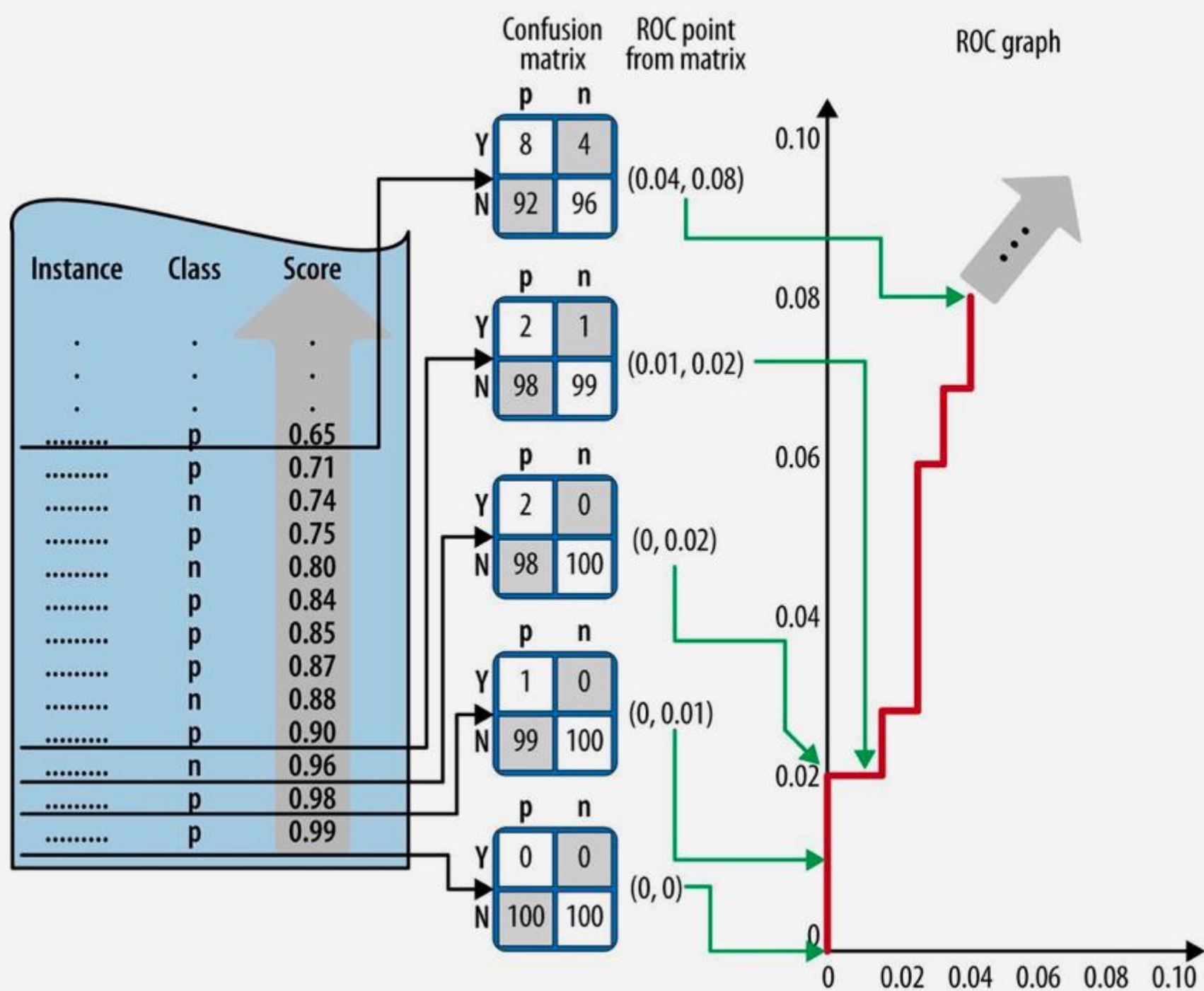
Benefit



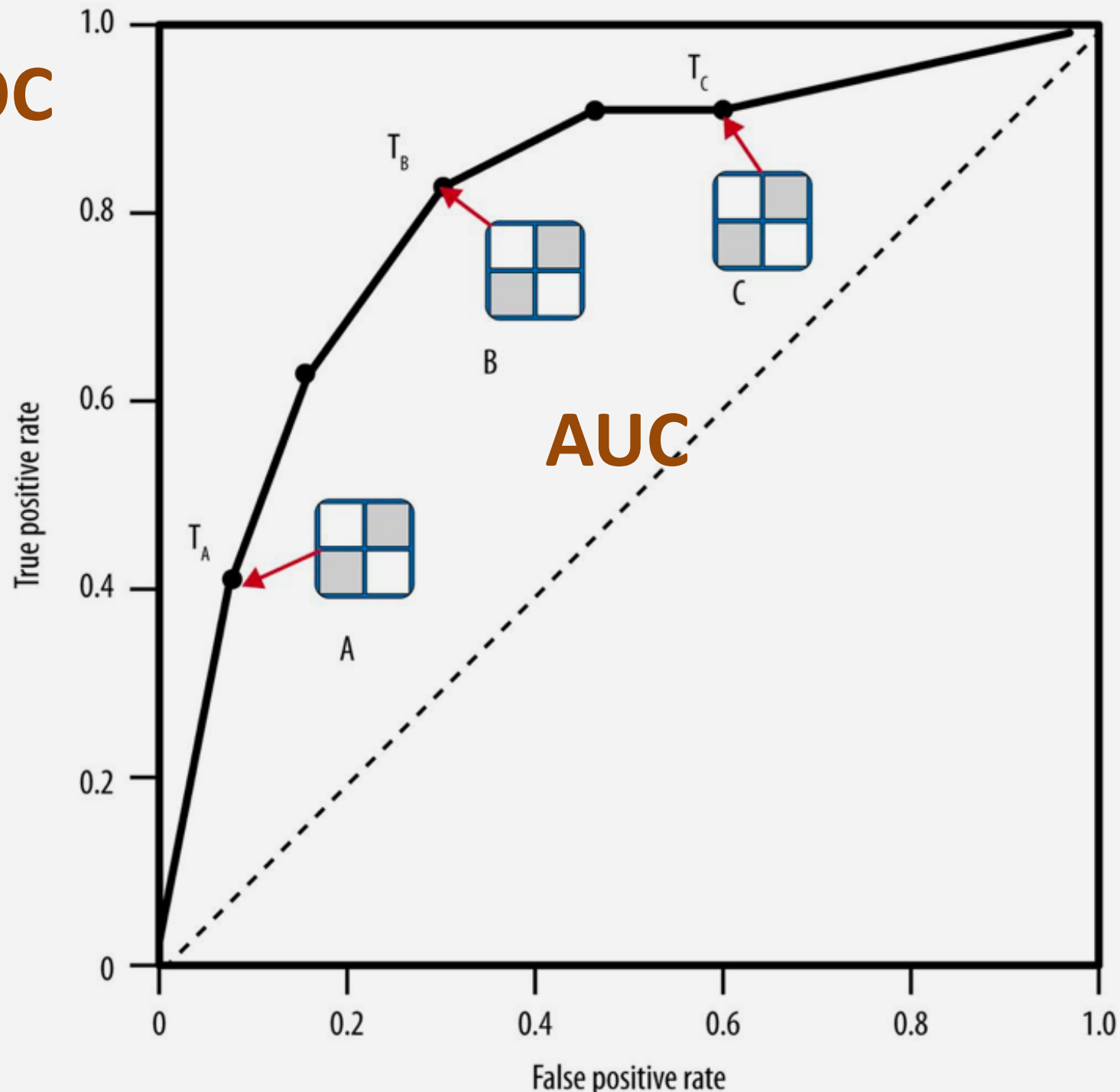
Cost

ROC

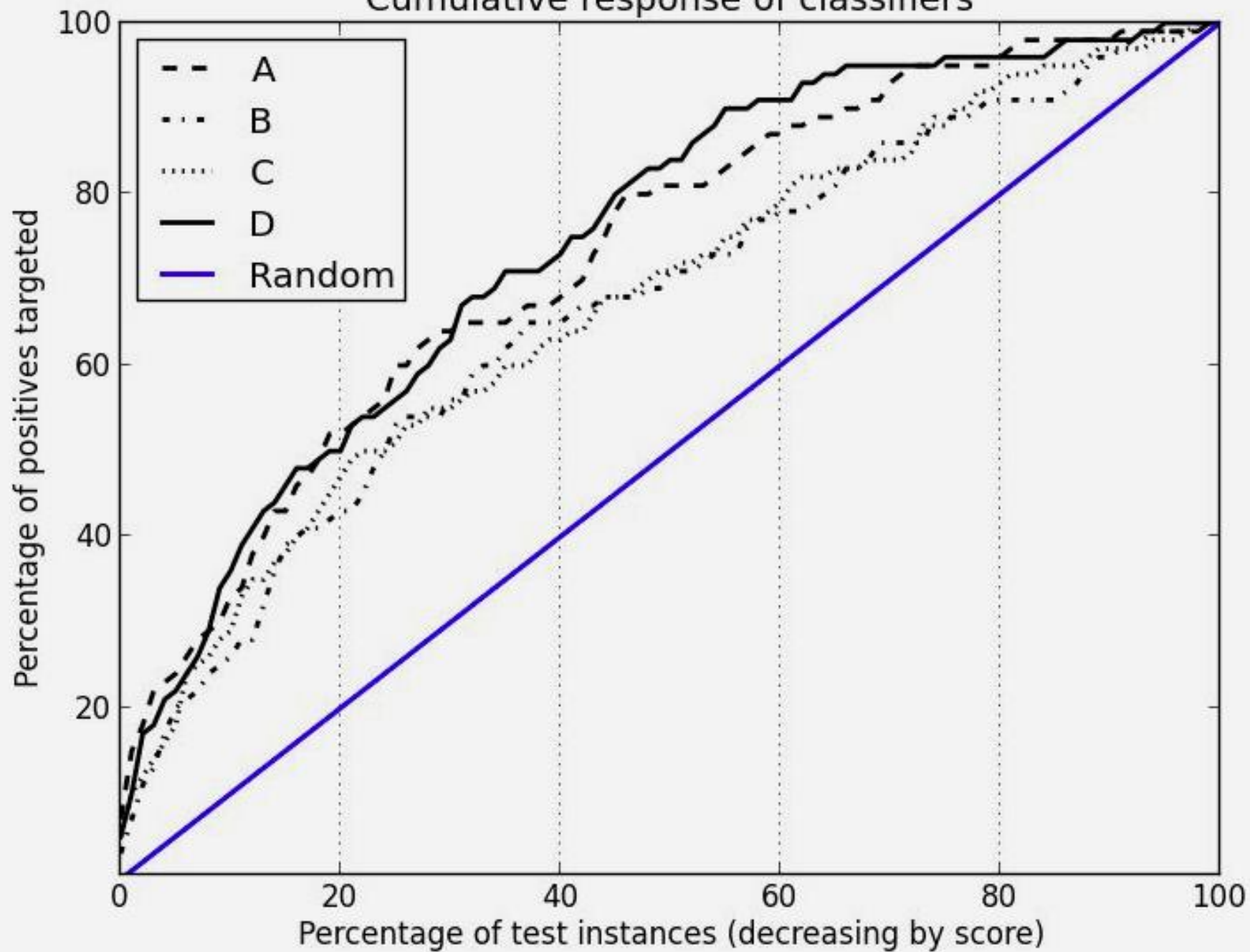
- $(0,0)$: never issuing a positive
- $(1,1)$: unconditionally issuing positives
- $(1,0)$: perfect classification
- $(0.5,0.5)$: guessing positive 50% of the time
- $(0.9,0.9)$: guessing positive 90% of the time
- Many real-world domains are dominated by negative cases, so the performance in the far left-hand side of the ROC graph is more interesting than elsewhere.



ROC



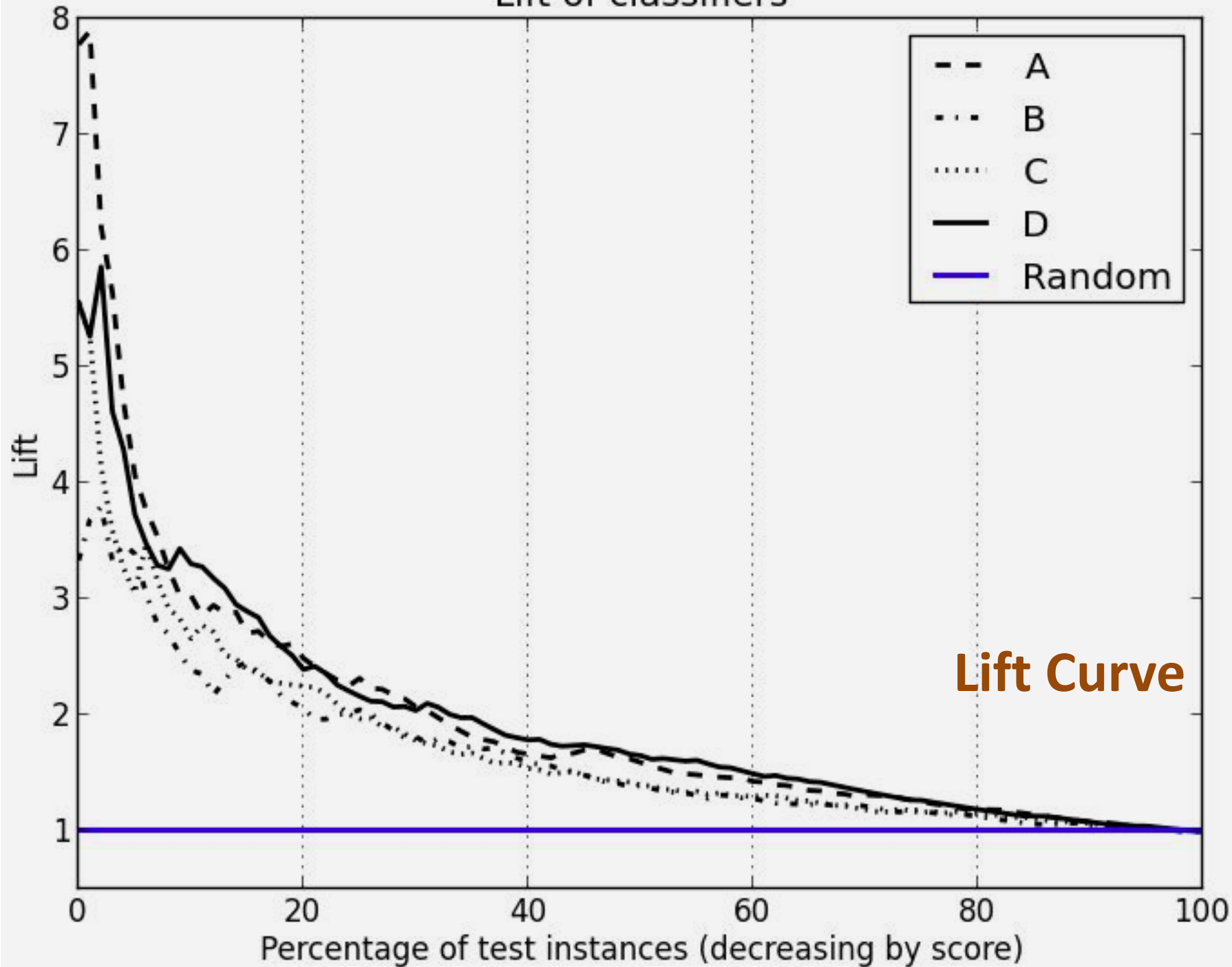
Cumulative response of classifiers



Lift Curve

- The lift of a classifier represents the advantage it provides over random guessing.
- Consider a list of 100 customers, half of them are positive instances. If you scan down the list and stop halfway.....
- If the list is sorted randomly, you would expect to see only half the positives. (a lift of 1).
- If the classifier is perfect, the lift is 2.
- Unlike for ROC curves, the lift curve assumes that the test set has exactly the same target class priors as the population to which the model will be applied.

Lift of classifiers



Business Problems and Analysis Tasks

Data-Analytic Thinking

- Business problems rarely are classification problems, regression problems or clustering problems.
- View business problems from a data perspective with structure and principles to guide you to the solutions.
- Bring together intuition, creativity, common sense and domain knowledge.

Targeting the best prospects for a charity mailing

- Would we want to model the probability of responding to the offer for each prospective donor.
- Would we like to maximize the total amount of donations? (the amount for this campaign or for the lifetime)
- Would we like to maximize the total amount of profit?

Expected benefit of targeting = $p(R \mid \mathbf{x}) \cdot v_R + [1 - p(R \mid \mathbf{x})] \cdot v_{NR}$

Expected benefit of targeting = $p(R \mid \mathbf{x}) \cdot v_R(\mathbf{x}) + [1 - p(R \mid \mathbf{x})] \cdot v_{NR}(\mathbf{x})$

$$\begin{aligned} p(R \mid \mathbf{x}) \cdot (d_R(\mathbf{x}) - c) + [1 - p(R \mid \mathbf{x})] \cdot (-c) &> 0 \\ p(R \mid \mathbf{x}) \cdot d_R(\mathbf{x}) - p(R \mid \mathbf{x}) \cdot c - c + p(R \mid \mathbf{x}) \cdot c &> 0 \\ p(R \mid \mathbf{x}) \cdot d_R(\mathbf{x}) &> c \end{aligned}$$

Selection Bias

- The data we are using to build the model may well be biased – meaning that they are selected randomly from the population to which you intend to apply the model.
- The data are from past donation – from the individual who did respond in the past.
- Some people may donate \$10 each and every time they're asked. Others may give \$100 and then feel they need not donate for a while.

Targeting the offer to customers who would churn when their contracts expire

- Do we really want to target our offer to those with the highest probability of defection?
- Targeting those with the highest value?
- Targeting those whom we would lose the most value if they were to leave.
- Assessing the expected benefit of not targeting.

$$EB_T(\mathbf{x}) = p(S \mid \mathbf{x}, T) \cdot (u_S(\mathbf{x}) - c) + [1 - p(S \mid \mathbf{x}, T)] \cdot (u_{NS}(\mathbf{x}) - c)$$

$$EB_{notT}(\mathbf{x}) = p(S \mid \mathbf{x}, notT) \cdot u_S(\mathbf{x}) + [1 - p(S \mid \mathbf{x}, notT)] \cdot u_{NS}(\mathbf{x})$$

$$\begin{aligned} VT &= p(S \mid \mathbf{x}, T) \cdot u_S(\mathbf{x}) - p(S \mid \mathbf{x}, notT) \cdot u_S(\mathbf{x}) - c \\ &= [p(S \mid \mathbf{x}, T) - p(S \mid \mathbf{x}, notT)] \cdot u_S(\mathbf{x}) - c \\ &= \Delta(p) \cdot u_S(\mathbf{x}) - c \end{aligned}$$

Competitive Advantage

Differences

- Netflix vs. Blockbuster
- Amazon vs. Barnes & Noble
- Amazon vs. eBay

Sustaining competitive advantages via data science

- Data and data science capabilities are assets to every company? (Web to Dell vs. Web to Compaq)
- Your competitors enjoy the same value from these assets?
- Formidable historical advantage; unique intellectual property; complementary assets

KSF

- Data and the capability to extract useful knowledge from data should be regarded as key strategic assets.
- Management must think data-analytically
- Management must create a culture to nurture data science and data scientists

The material is excerpted from
“Data Science for Business”
Foster Provost and Tom Fawcett