


Business Intelligence

13 What is a good model?

Prof. Dr. Bastian Amberg
(summer term 2024)

28.6.2024

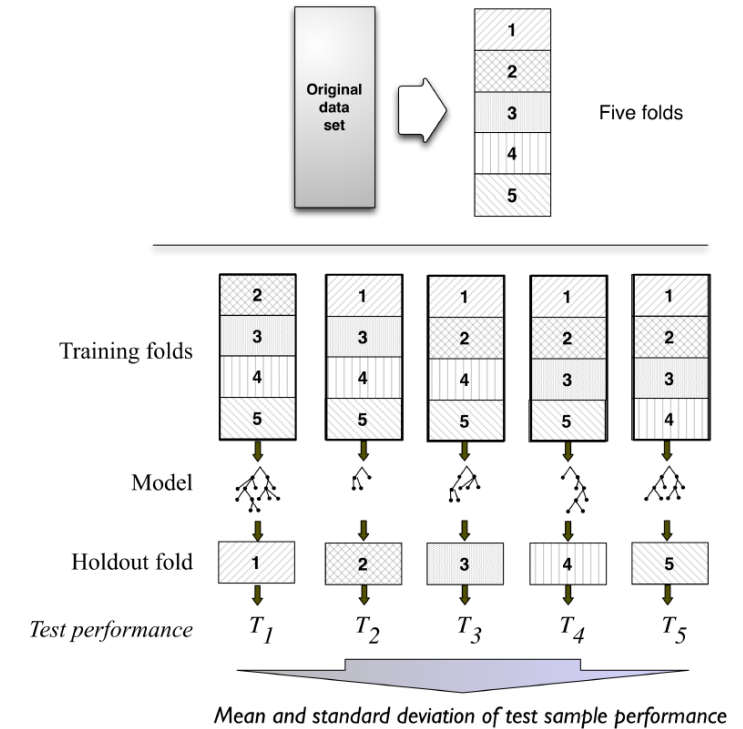
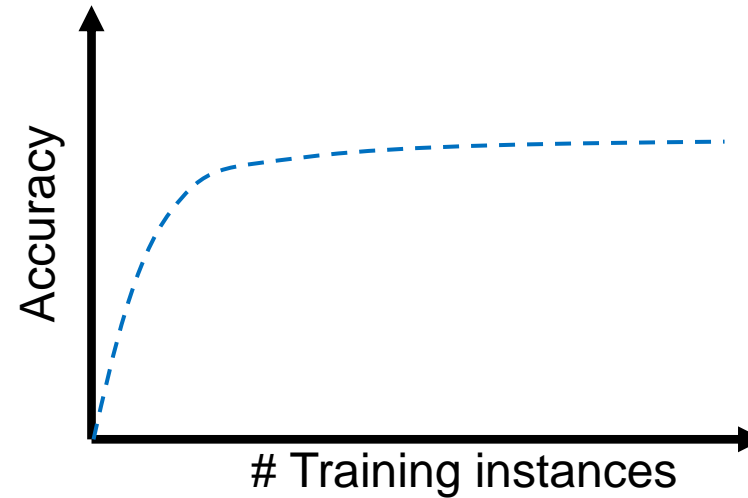
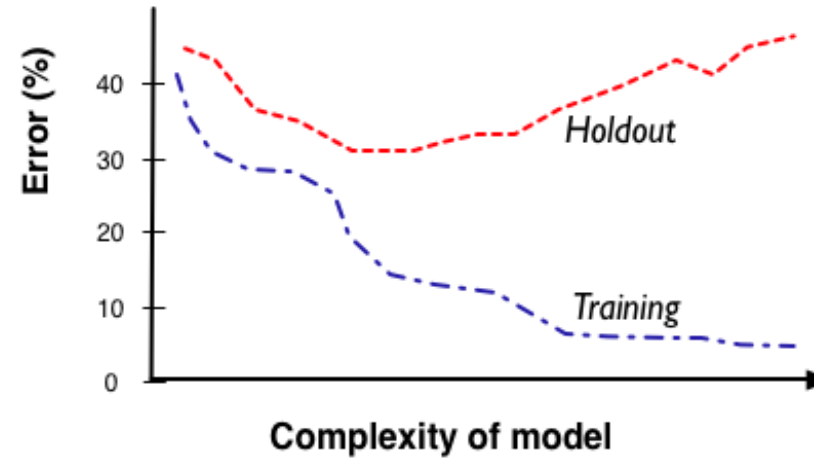
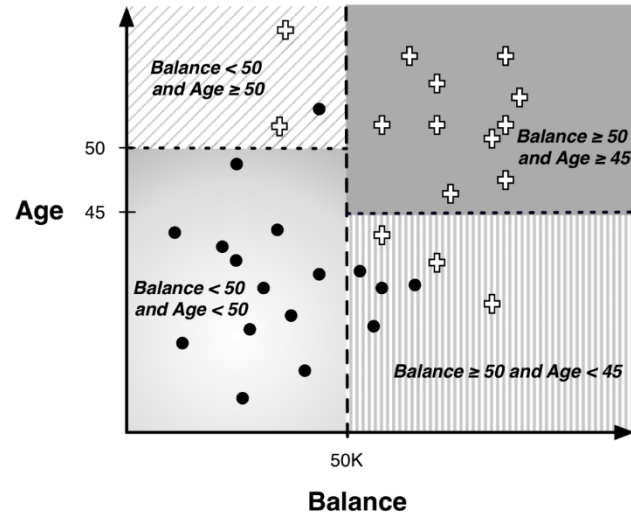
Schedule

		Wed., 10:00-12:00			Fr., 14:00-16:00 (Start at 14:30)		Self-study	
Basics	W1	17.4.	(Meta-)Introduction		19.4.		Python-Basics	Chap. 1
	W2	24.4.	Data Warehouse – Overview	& OLAP	26.4.	[Blockveranstaltung SE Prof. Gersch]		Chap. 2
	W3	1.5.			3.5.			Chap. 3
	W4	8.5.	Data Warehouse Modeling I	& II	10.5.	Data Mining Introduction		
Main Part	W5	15.5.	CRISP-DM, Project understanding		17.5.	Python-Basics-Online Exercise	Python-Analytics	Chap. 1
	W6	22.5.	Data Understanding, Data Visualization I		24.5.	No lectures, but bonus tasks 1.) Co-Create your exam 2.) Earn bonus points for the exam		Chap. 2
	W7	29.5.	Data Visualization II		31.5.			
	W8	5.6.	Data Preparation		7.6.	Predictive Modeling I (10:00 -12:00)	BI-Project	Start
	W9	12.6.	Predictive Modeling II		14.6.	Python-Analytics-Online Exercise		
	W10	19.6.	Guest Lecture Dr. Ionescu		21.6.	Fitting a Model		
	W11	26.6.	How to avoid overfitting		28.6.	What is a good Model?		
Deepening	W12	3.7.	Project status update Evidence and Probabilities		5.7.	Similarity (and Clusters) From Machine to Deep Learning I		
	W13	10.7.			12.7.	From Machine to Deep Learning II		
	W14	17.7.	Project presentation		19.7.	Project presentation		End
Ref.						Klausur 1. Termin, 31.7. '24 Klausur 2. Termin, 2.10. '24	Projektbericht	

Case Study

Last Lesson(s) – Exercise

Explain what you see.



Kahoot-Fragen
www.kahoot.it
 (über Smartphone oder Laptop)
 PIN folgt

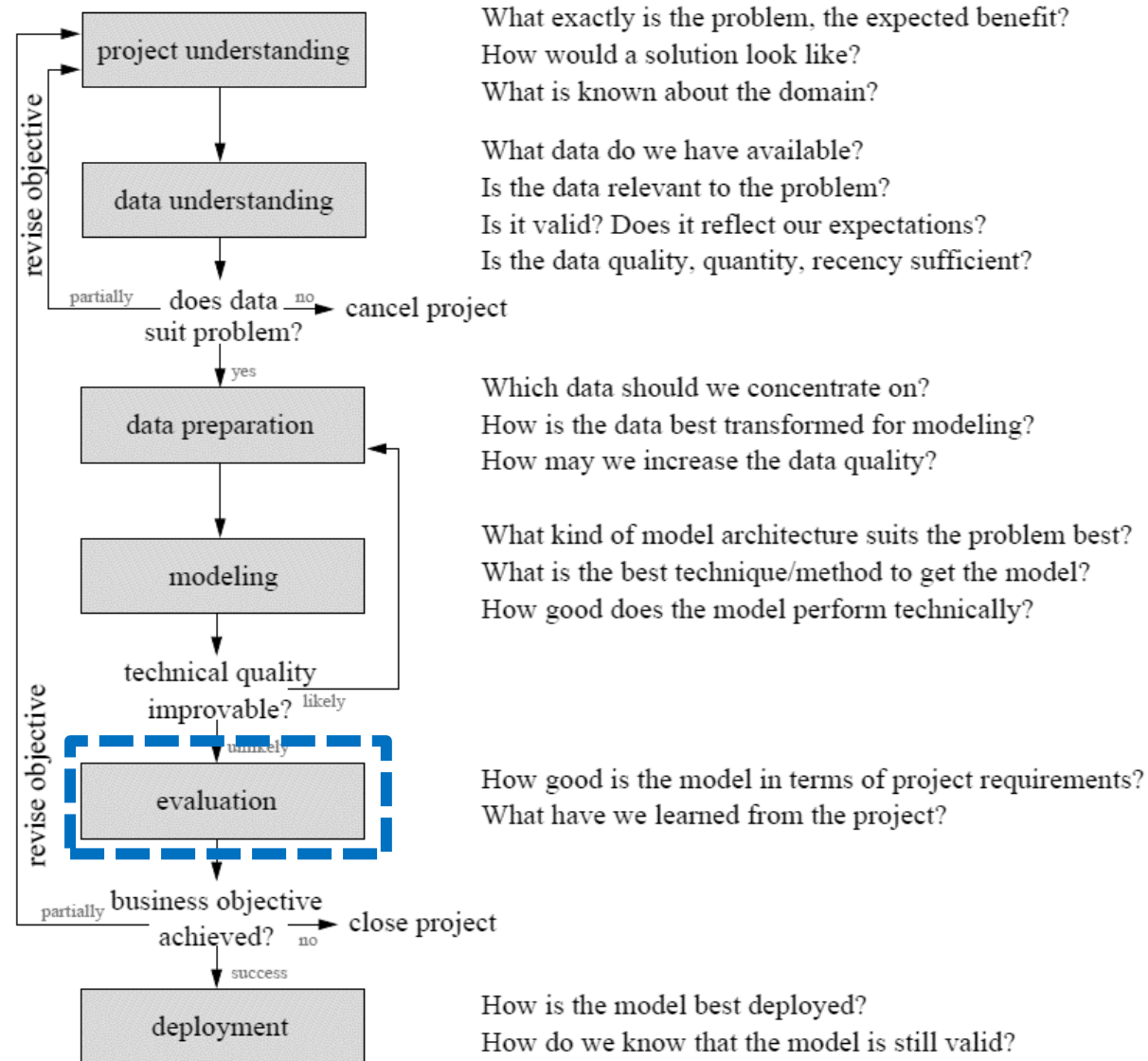
10 Min.

Cross
Industry
Standard
Process for
Data
Mining

Iteration as
a rule

Process of data
exploration

Implementation of the
KDD Process



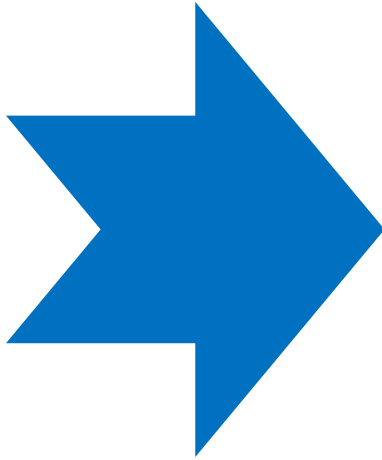
Modeling: Open issue from last lesson:

Avoiding Overfitting and Complexity Control
in Parametric Learning/Optimization

$$\arg \max_{\mathbf{w}} \text{fit}(\mathbf{x}, \mathbf{w}) \rightarrow \text{penalize complexity}$$

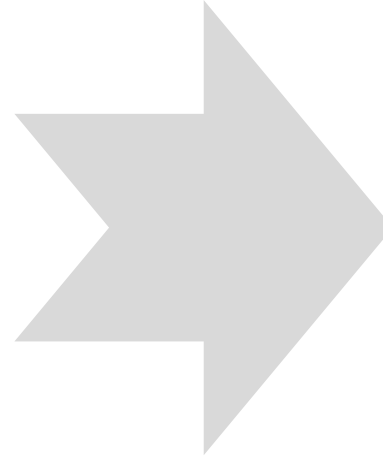
$$\arg \max_{\mathbf{w}} [\text{fit}(\mathbf{x}, \mathbf{w}) - \lambda \cdot \text{penalty}(\mathbf{w})]$$

we will come back to this topic later



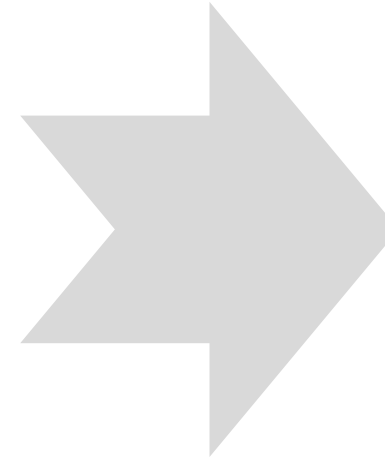
(1) Measuring accuracy

- Confusion matrix
- Unbalanced classes



(2) Expected Value

- Evaluate classifier use
- Frame classifier evaluation



(3) Evaluation and baseline performance

What is desired from data mining results?

How would you **measure** that your model is any good?
How to measure performance in a meaningful way?

Model evaluation is **application-specific**

We look at common issues and themes in evaluation

Frameworks and metrics for classification and instance scoring

Think about the specific BI project you are working on....



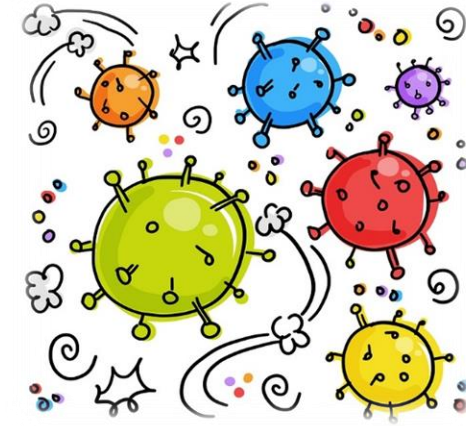
Bad positives and harmless negatives

Classification terminology

a bad outcome	→ a “positive” example	[alarm!]
a good outcome	→ a “negative” example	[uninteresting]

Further examples

medical test:	positive test → disease is present
fraud detector:	positive test → unusual activity on account



A classifier tries to distinguish the majority of cases (**negatives**, the uninteresting) from the small number of alarming cases (**positives**, alarming)

number of mistakes made on **negative** examples (false positive errors) will be relatively high / may dominate

cost of each mistake made on a **positive** example (false negative error) will be relatively high / will be higher

Measuring accuracy and its problems

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

The 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 labels with predicted classes.

For a binary variable it is as follows:

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

Each example in a test set has an **actual class label** and the **class predicted** by the classifier

Interpretation:

The confusion matrix separates the decisions made by the classifier

- **actual/true classes:** p(ositive), n(egative)
- **predicted classes:** Y(es), N(o)
- The main diagonal contains the count of correct decisions

Mini-Exercise - The confusion matrix

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

Diese Folie ist nach der Vorlesung mit Lösungen verfügbar

Kahoot-Fragen
www.kahoot.it
(über Smartphone oder Laptop)
PIN folgt

10 Min.

The confusion matrix

Unbalanced classes

In most real world classification problems, one class is often **rare**

Classification is used to find a relatively small number of **unusual ones** (defrauded customers, defective parts, targeting consumers who actually would respond, ...)

The class distribution is unbalanced (“skewed”)

Evaluation based on **accuracy** does not work

Example: 999:1 ratio
always choose the most prevalent class – 99.9% accuracy!

Fraud detection: skews of 1:99
Is a model with 80% accuracy always better than a model with 37% accuracy?

We need to have more information about the population



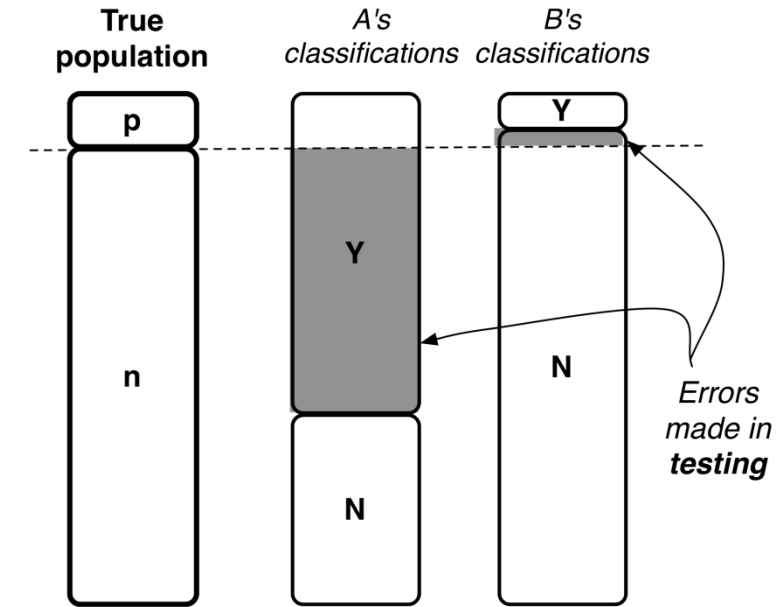
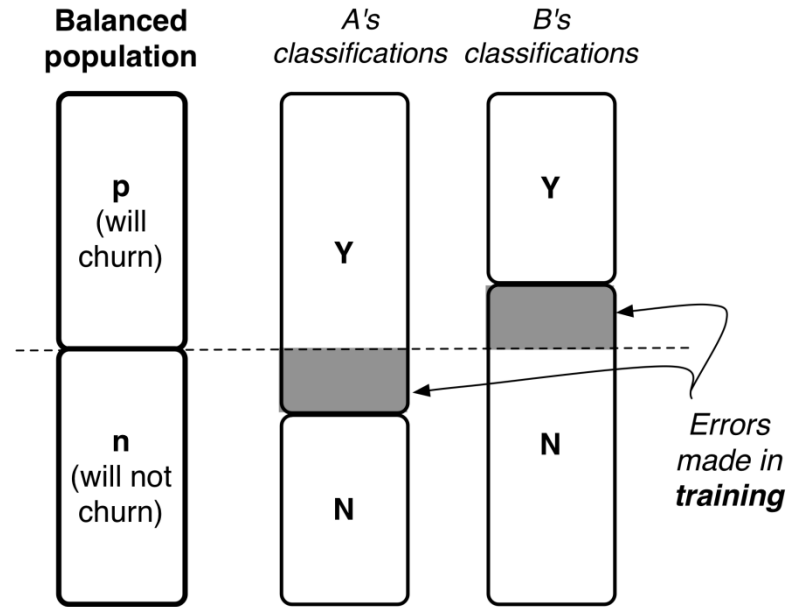
The confusion matrix

Unbalanced classes

Consider two models A and B for the churn example (~every tenth customer churns)

We train with a balanced population (1000 customer)

- Both models correctly classify 80% of the balanced pop.
- Classifier A often falsely predicts that customers will churn
- Classifier B makes many opposite errors



Unbalanced population:

- A's accuracy is 40%,
B's accuracy is 97%

Note the **different performances** of the models in form of a confusion matrix:

$$CM_A = \begin{matrix} & \text{churn} & \text{not churn} \\ \begin{matrix} Y \\ N \end{matrix} & \begin{pmatrix} 500 & 200 \\ 0 & 300 \end{pmatrix} \end{matrix}$$

$$CM_B = \begin{matrix} & \text{churn} & \text{not churn} \\ \begin{matrix} Y \\ N \end{matrix} & \begin{pmatrix} 300 & 0 \\ 200 & 500 \end{pmatrix} \end{matrix}$$

$$CM_A = \begin{matrix} & \text{churn} & \text{not churn} \\ \begin{matrix} Y \\ N \end{matrix} & \begin{pmatrix} 100 & 600 \\ 0 & 300 \end{pmatrix} \end{matrix}$$

$$CM_B = \begin{matrix} & \text{churn} & \text{not churn} \\ \begin{matrix} Y \\ N \end{matrix} & \begin{pmatrix} 70 & 0 \\ 30 & 900 \end{pmatrix} \end{matrix}$$

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 she/he does not)

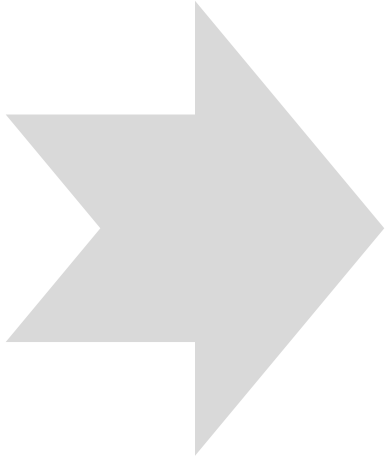
→ **false positive error**, expensive, but not life threatening

a patient has cancer, but she/he is told that she/he has not

→ **false negative error**, more serious

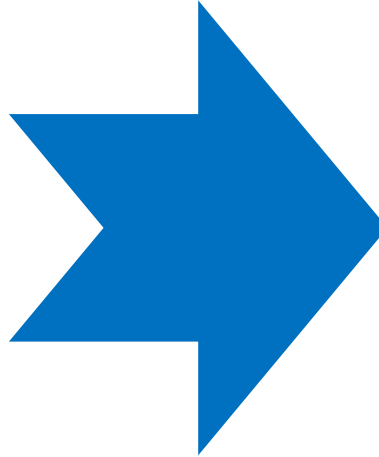
Errors should be counted separately:

Estimate cost or benefit of each decision



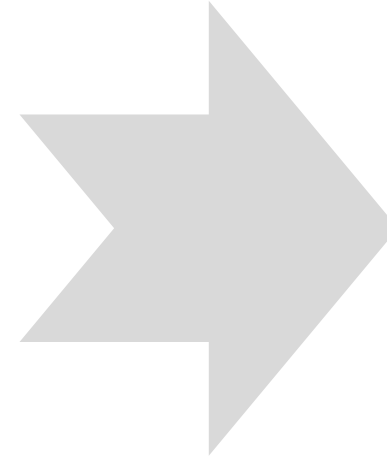
(1) Measuring accuracy

- Confusion matrix
- Unbalanced classes



(2) Expected Value

- Evaluate classifier use
- Frame classifier evaluation



(3) Evaluation and baseline performance

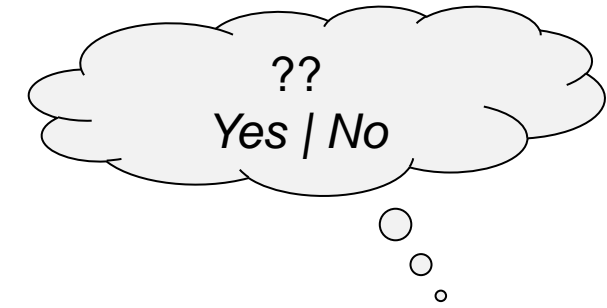
The expected value framework

Expected value calculation includes **enumeration of the possible outcomes** of a situation

Expected value = weighted average of the values of different possible outcomes, where the weight given to each value is the probability of its occurrence

Example: different levels of profit

We focus on the maximization of expected profit



General form of expected value computation:

$$EV = p(o_1) \cdot v(o_1) + p(o_2) \cdot v(o_2) + \dots +$$

with o_i as possible decision outcome,
 $p(o_i)$ as its probability, and $v(o_i)$ as its value.

Probabilities can be **estimated** from available data

- We consider two cases
- I. Evaluate classifier use
 - II. Frame classifier evaluation

(I) Expected value for use of a classifier

Use of a classifier: predict a class and take some action

Example target marketing: assign each consumer to either a class „likely responder“ or „not likely responder“

Response is usually relatively low – so no consumer may seem like a likely responder

Computation of the expected value

A model gives an estimated probability of response $\hat{p}_R(x)$ for any consumer with a feature vector x

Calculate expected benefit (or costs) of targeting consumer x : $\hat{p}_R(x) \cdot v_R + (1 - \hat{p}_R(x)) \cdot v_{NR}$ with v_R being the value of a response and v_{NR} the value from no response

Example:

Price of product: \$200, costs of product: \$100

Targeting a consumer: \$1, profit $v_R = \$99$, $v_{NR} = -\$1$

Do we make a profit? Is the expected value (profit) of targeting greater than zero?

$$\hat{p}_R(x) \cdot \$99 + (1 - \hat{p}_R(x)) \cdot (-\$1) > 0$$

$$\Leftrightarrow \hat{p}_R(x) \cdot \$99 > (1 - \hat{p}_R(x)) \cdot \$1$$

$$\Leftrightarrow \hat{p}_R(x) > 0.01$$

We should target the consumer as long as the estimated probability of responding is greater than 1%

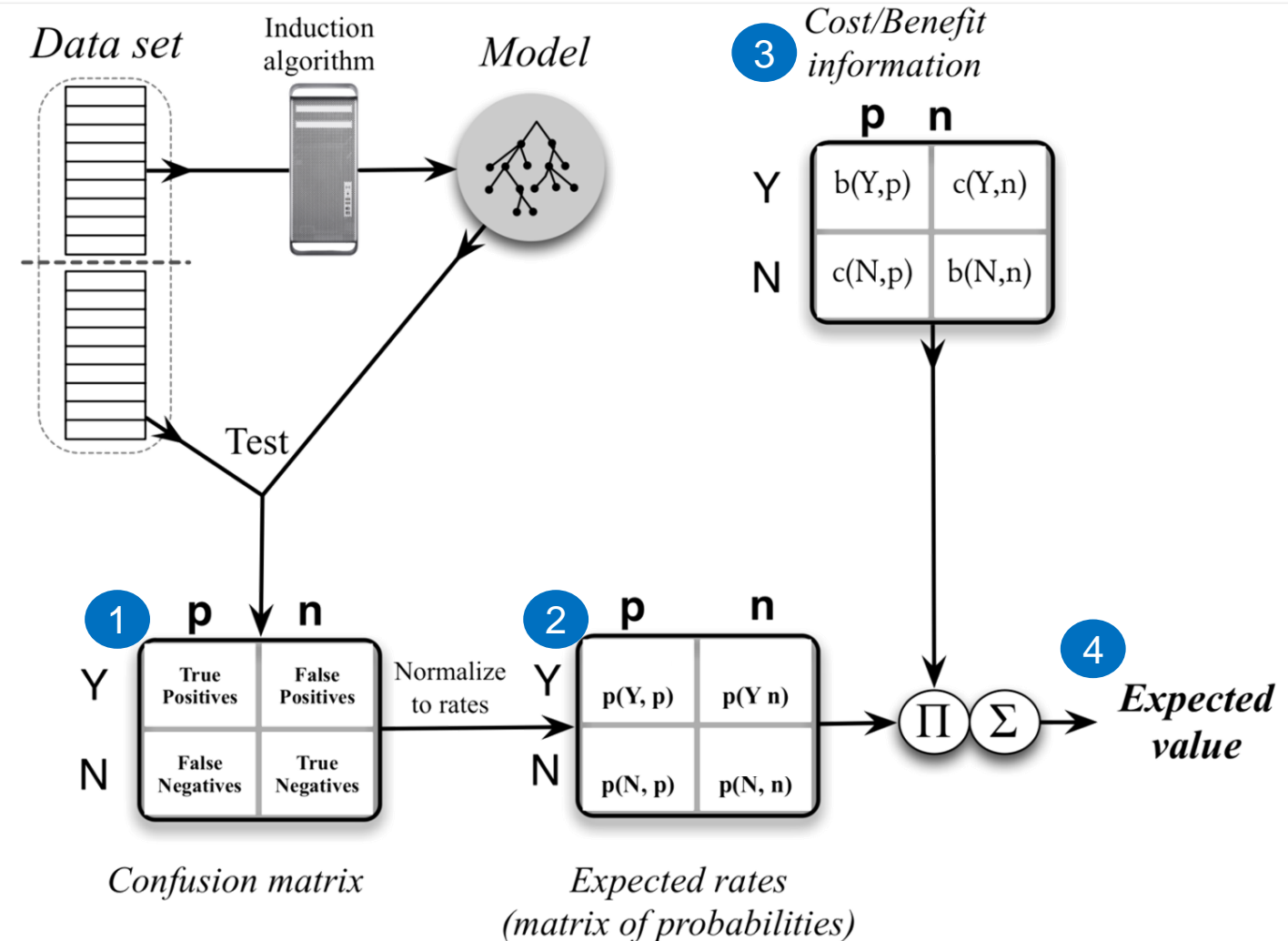
(II) Expected value classification

Goal: **compare the quality of different models** with each other

- Does the data-driven model perform better than a hand-crafted model?
- Does a classification tree work better than a linear discriminant model?
- Do any of the models perform substantially better than a baseline model?

In short:

How well does each model perform with regards to its expected value?



Expected rates for evaluation of a classifier

2

Aggregate together all the different cases:

When we target consumers, what is the probability that they (do not) respond?

What about when we do not target consumers, would they have responded?

This information is available in the **confusion matrix**

Each o_i corresponds to one of the possible combinations of the class we predict/the actual class

Where do the probabilities of errors and correct decisions actually come from?

Each cell of the confusion matrix contains a count of the number of decisions corresponding to the combination of (predicted, actual)

$count(h, a)$

Compute estimated probabilities as

$$p(h, a) = count(h, a) / Total$$

Example confusion matrix/estimates of probability

Predicted	Actual	
	p	n
	Y	N
Y	56	7
N	5	42

$T = 110, P = 61, N = 49$ (**P**ositive, **N**egative)

$$p(Y, p) = \frac{56}{110} = 0.51, \quad p(Y, n) = \frac{7}{110} = 0.06$$

$$p(N, p) = \frac{5}{110} = 0.05, \quad p(N, n) = \frac{42}{110} = 0.38$$

Costs and benefits

3

Compute **cost-benefit values** for each decision pair

A cost-benefit matrix specifies for each (predicted,actual) pair the cost or benefit for making such a decision

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

Correct classifications (true positives and negatives) correspond to $b(Y,p)$ and $b(N,n)$, respectively

Incorrect classifications (false positives and negatives) correspond to $c(Y,n)$ and $c(N,p)$, respectively
[often negative benefits or costs]

Costs and benefits cannot be estimated from data

How much is the true value for retaining a customer? (i.e., CLV)

Ref. Often use of average estimated costs and benefits

Targeted marketing example

- **False positive** occurs when we classify a consumer as a likely responder and therefore target her, but she does not respond
→ cost $c(Y,n) = -1$ //or negative benefit $b(Y,n)$
- **False negative** is a consumer who was predicted not to be a likely responder, but would have bought if offered. No money spent, nothing gained
→ cost $c(N,p) = 0$ //or negative benefit $b(N,p)$
- **True positive** is a consumer who is offered the product and buys it
→ benefit $b(Y,p) = 200 - 100 - 1 = 99$
- **True negative** is a consumer who was not offered a deal but who would not have bought it
→ benefit $b(N,n) = 0$

Sum up in cost-benefit matrix

Predicted	Actual	
	p	n
	Y	N
Y	99	-1
N	0	0

Expected profit computation

4

Sufficient for comparison of various models

Compute **expected profit** by cell-wise multiplication of the matrix of costs and benefits against the matrix of probabilities:

$$EP = p(Y, p) \cdot b(Y, p) + p(N, p) \cdot c(N, p) + \\ p(N, n) \cdot b(N, n) + p(Y, n) \cdot c(Y, n)$$

	p	n
Y	$p(Y, p)$	$p(Y, n)$
N	$p(N, p)$	$p(N, n)$

	p	n
Y	$b(Y, p)$	$c(Y, n)$
N	$c(N, p)$	$b(N, n)$

Alternative calculation: factor out the probabilities of seeing each class (class priors) [Use $p(x, y) = p(y) \cdot p(x | y)$]

Class priors $p(p)$ and $p(n)$ specify the likelihood of seeing positive versus negative instances

Factoring out allows us to separate the influence of class imbalance from the predictive power of the model

Factoring out priors yields the following

alternative expression for expected profit:

$$EP = p(Y|p) \cdot p(p) \cdot b(Y, p) + p(N|p) \cdot p(p) \cdot c(N, p) + \\ p(N|n) \cdot p(n) \cdot b(N, n) + p(Y|n) \cdot p(n) \cdot c(Y, n)$$

$$EP = p(p) \cdot [p(Y|p) \cdot b(Y, p) + p(N|p) \cdot c(N, p)] + \\ p(n) \cdot [p(N|n) \cdot b(N, n) + p(Y|n) \cdot c(Y, n)]$$

The first component corresponds to the expected profit from the **positive examples**, whereas the second corresponds to the expected profit from the **negative examples**

We call:

$p(Y|p)$: true positive rate

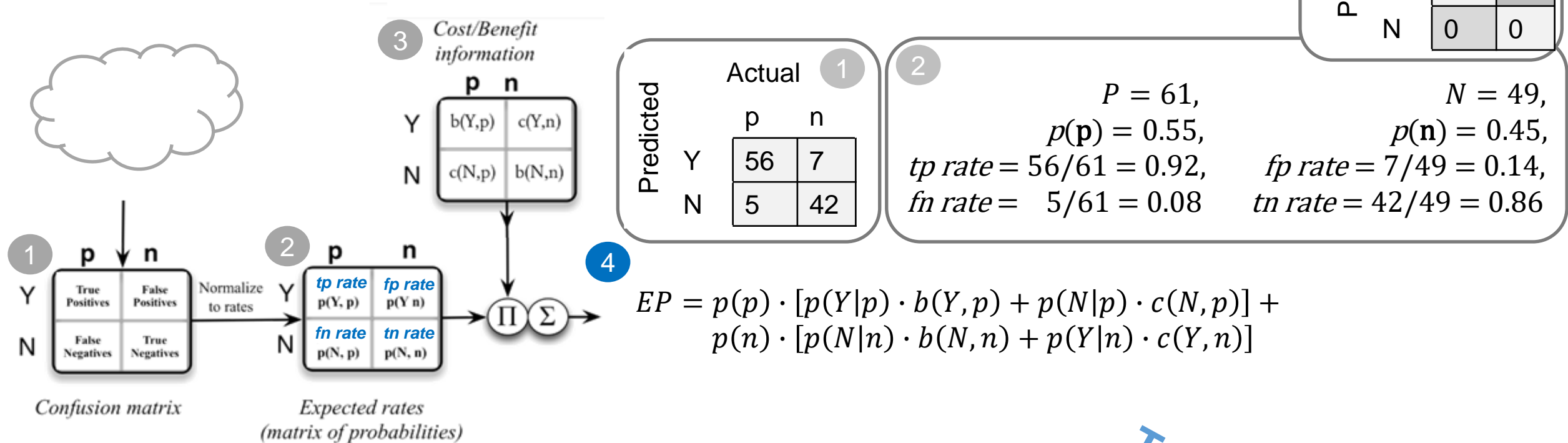
$p(N|p)$: false negative rate

$p(Y|n)$: false positive rate

$p(N|n)$: true negative rate

Exercise – Expected value computation

Example alternative expression



		Actual 1	
		p	n
Predicted	Y	56	7
	N	5	42

		Actual 2	
		p	n
Predicted	Y	99	-1
	N	0	0

		Actual 3	
		p	n
Predicted	Y	99	-1
	N	0	0

2

$P = 61,$
 $p(p) = 0.55,$
 $tp\ rate = 56/61 = 0.92,$
 $fn\ rate = 5/61 = 0.08$

$N = 49,$
 $p(n) = 0.45,$
 $fp\ rate = 7/49 = 0.14,$
 $tn\ rate = 42/49 = 0.86$

This expected value means that ...

Todo for Wednesday

5 Min.

Based on the entries of the confusion matrix, we can describe various evaluation metrics

Accuracy (count of correct decisions): $\frac{TP+TN}{P+N}$

True positive rate / Recall / Specificity : $\frac{TP}{TP+FN}$

False negative rate: $\frac{FN}{TP+FN}$

Sensitivity: $\frac{TN}{TN+FP}$

Precision (accuracy over the cases predicted to be positive): $\frac{TP}{TP+FP}$

F-measure (harmonic mean): $2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$

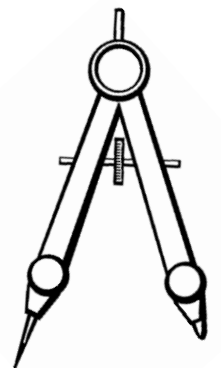
		p	n
Predicted Y		True positives	False positives
N		False negatives	True negatives

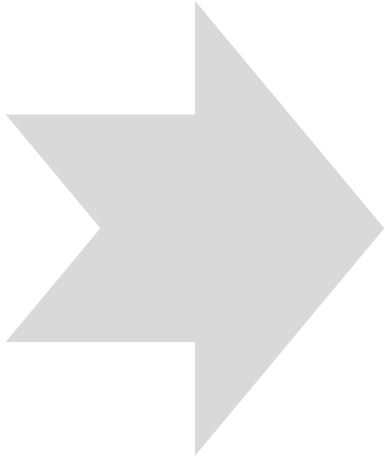
		p	n
Predicted Y		True positives	False positives
N		False negatives	True negatives

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

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

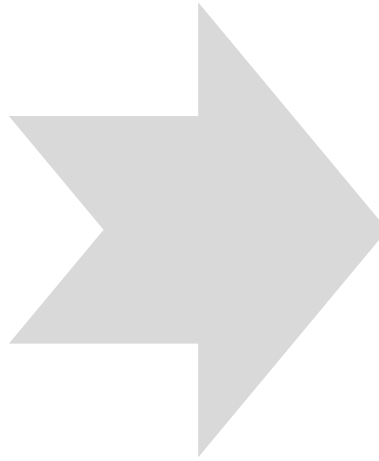
		p	n
Predicted Y		True positives	False positives
N		False negatives	True negatives





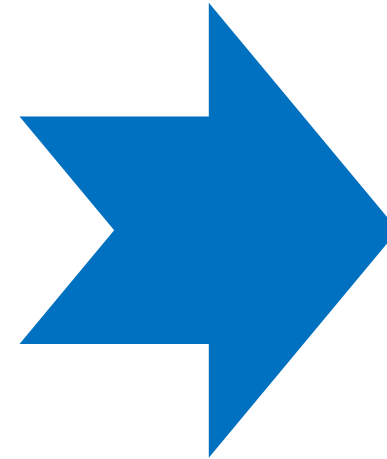
(1) Measuring accuracy

- Confusion matrix
- Unbalanced classes



(2) Expected Value

- Evaluate classifier use
- Frame classifier evaluation



(3) Evaluation and baseline performance

Baseline performance (1/3)

Consider what would be a **reasonable baseline** against which to compare model performance

Demonstrate stakeholder that data mining has added value (or not)

What is the appropriate baseline for comparison?

Depends on the actual application

There are two basic tests that any weather forecast must pass to demonstrate its merit:

*(1) It must do **better than** what meteorologists call **persistence**: the assumption that the weather will be the same tomorrow (and the next day) as it was today.*

*(2) It must also **beat** climatology, the **long-term historical average** of conditions on a particular date in a particular area (not only dependent to time/seasonal effects).*



Baseline performance (2/3)

Baseline performance for classification

Compare to a completely random model (very easy)

Implement a simple (but not simplistic) alternative model

Majority classifier = a naive classifier that always chooses the majority class of the training data set

May be challenging to outperform: classification accuracy of 94%, but only 6% of the instances are positive

→ majority classifier also would have an accuracy of 94%!

Pitfall: don't be surprised that many models simply predict everything to be of the **majority class**

Maximizing simple prediction accuracy is usually not an appropriate goal



Hint:



DummyClassifier is a classifier that makes predictions using simple rules. This classifier is useful as a simple baseline to compare with other (real) classifiers.

Strategies, e.g.,

“most_frequent”:

always predicts the most frequent label in the training set.

“uniform”:

generates predictions uniformly at random.

(<https://scikit-learn.org/stable/modules/generated/sklearn.dummy.DummyClassifier.html>)

Further alternative:

how well does a simple “conditional” model perform?

Conditional → prediction different based on the value of the features

Just **use the most informative variable** for prediction

Decision tree: build a tree with only one internal node (decision stump) → tree induction selects the single most informative feature to make a decision

Compare quality of models based on data sources

Quantify the value of each source

Implement models that are based on **domain knowledge**

Fragen?

- ✓ Measuring accuracy
 - ✓ Confusion matrix
 - ✓ Unbalanced classes
- ✓ A key analytical framework: Expected value
 - ✓ Evaluate classifier use
 - ✓ Frame classifier evaluation
- ✓ Evaluation and baseline performance

- Please remember:
*Gemäß Vorlesungsplanung ist für nächsten Mittwoch **Project status update** vorgesehen.*
Was heißt das nun eigentlich?
 - Bitte jede Projektgruppe **bis Mittwoch 3.7.** einen kurzen „Zwischenbericht“ liefern
 - Formlos per E-Mail an bastian.amberg@fu-berlin.de , Betreff „BI Projektgruppe – Zwischenstand“
 - Aktueller Stand (z.B. durchgeführte Bearbeitungsschritte nach CRISP-DM)
 - Offene Punkte bzw. grobes weiteres Vorgehen
 - **kurz (d.h. ca. 4-5 Sätze)**
 - Optionaler weiterer Punkt: Besteht Sprechstundenbedarf?
Sprechstunde dann in der Woche vom 8.7. bis 12.7.
- Please do the “Exercise – Expected value computation”
Slide 23

Recommended reading

What is a good model:

- | | |
|-----------------|---|
| Provost, F., | Data Science for Business |
| Fawcett, T. | Chapter 7 |
| Berthold et al. | Guide to Intelligent Data Analysis, Chapter 5 |

Further reading (beyond this class):

- | | |
|--------------|---|
| Provost, F., | Data Science for Business |
| Fawcett, T. | Chapter 8, Visualizing Model Performance (discusses graphical views of model behaviour) |



Business Intelligence – Interim summary

We have three goals. After this course:

- You know how to solve business problems by **data-analytic thinking**
- You know **tools** and ways of how to practically **implement** solution methods
- You have an overview about principles of **how to model and solve** upcoming **business problems**.

Main focus:

Data Warehousing / Data Engineering ► How to **store and access** huge amounts of data?

- DW Overview (incl. OLAP) L02
- DW Modelling L03, L04

Data Mining / Data Science ► How to **derive knowledge and profitable business action** out of (large) databases?

- DM Overview L04, L05
- CRISP-DM

- project understanding L05
- data understanding L06, L07
- data preparation L08
- modeling L09, L10, L11, L12, ...
- evaluation L13 ...

Exkurs: Weiterführend siehe Schulz et al. 2020 DASC-PM

