
Web-Mining Agents

Prof. Dr. Ralf Möller

Universität zu Lübeck

Institut für Informationssysteme

Karsten Martiny (Übungen)



Organizational Issues: Assignments

- **Start:** Wed, 21.10., 2-4pm, AMHZ S1, Class also Thu 2-4pm, IFIS 2035
- **Lab:** Fr. 2-4pm, Building 64, Inst. Math., Seminar room Hilbert (3rd floor) (registration via Moodle right after this class)
- **Assignments** provided via Moodle after class on Thu.
- **Submission of solutions** by Wed 2pm, small kitchen IFIS (one week after provision of assignments)
- **Work on assignments** can/should be done in groups of 2 (pls. indicate name and group on submitted solution sheets)
- In **lab classes on Friday**, we discuss assignments from current week and understand solutions for assignments from previous week(s)

Organizational Issues: Exam

- **Registration** in class required to be able to participate in **oral exam** at the end of the semester (2 slots)
- **Prerequisite** to participate in exam:
50% of all points of the assignments

Search Engines: State of the Art

- **Input:** Strings (typed or via audio), images, ...
- **Public services:**
 - Links to web pages plus mini synopses via GUI
 - Presentations of structured information via GUI excerpts from the Knowledge Vault
http://videolectures.net/kdd2014_murphy_knowledge_vault/
(previously known as Knowledge Graph)
- **NSA services: ?**
- **Methods:** Information retrieval, machine learning
- **Data:** Grabbed from free resources (win-win suggested)

Search Results

Web

Images

Maps

Videos

News

Shopping

More

Indiana, PA

Change location

Show search tools

<https://www.google.com/#hl=en>

pittsburgh - Google Search

File Edit View Favorites Tools Help

AVG Search Safe Do Not Track Weather Facebook Speedtest

[City of Pittsburgh, Pennsylvania - Pghgov.com](http://www.cityofpittsburgh.pa.us/)
www.cityofpittsburgh.pa.us/
Official city site including information on economic development, resident information, links, tourism and contact information.

[Pittsburgh Hotels, Attractions & Vacation Packages : Pittsburgh PA ...](http://www.visitpittsburgh.com/)
www.visitpittsburgh.com/
Greater Pittsburgh Convention & Visitors Bureau features vacation planning and travel information for Pittsburgh hotels, tours, attractions, meetings, restaurants, ...
[Things to Do - Get A Visitors Guide - Contact Us - Events](#)

[Pittsburgh - Wikipedia, the free encyclopedia](http://en.wikipedia.org/wiki/Pittsburgh)
en.wikipedia.org/wiki/Pittsburgh
Pittsburgh is the second-largest city in the U.S. Commonwealth of Pennsylvania and the county seat of Allegheny County. Regionally, it anchors the largest ...
[Port Authority of Allegheny ... - History - Neighborhoods - University of Pittsburgh](#)


[About Pittsburgh - Pittsburgh Pennsylvania for Residents & Visitors](http://pittsburgh.about.com/)
pittsburgh.about.com/
5 days ago - Get the scoop on everything Pittsburgh has to offer, from things to do and see, places to live or visit, upcoming events, unique attractions, ...

[Frommer's Pittsburgh](http://www.frommers.com)
www.frommers.com > ... > North America > USA > Pennsylvania
Plan your Pittsburgh vacation with the Frommer's comprehensive Pittsburgh travel guide. The travel guide includes information about hotels, restaurants, ...

[University of Pittsburgh](http://www.pitt.edu/)
www.pitt.edu/
The University of Pittsburgh is among the nation's most distinguished comprehensive universities, with a wide variety of high-quality programs in both the arts ...

[News for pittsburgh](#)

[Pittsburgh Pirates - TeamReport](#)
[Chicago Tribune](#) - 15 minutes ago
MLB Team Report - Pittsburgh Pirates - INSIDE PITCH.



Pittsburgh

Pittsburgh is the second-largest city in the U.S. Commonwealth of Pennsylvania and the county seat of Allegheny County. Regionally, it anchors the largest urban area of both Appalachia and the Ohio River Valley. Wikipedia

Founded: November 25, 1758

Area: 58.3 sq miles (151 km²)

Weather: 61°F (16°C), Wind S at 10 mph (16 km/h), 70% Humidity


Local time: 11:26pm Tuesday (EDT)

Population: 307,484 (2011)

Upcoming events

Sep 26 Wed	Pittsburgh Penguins vs. Detroit Red Wings CONSOL Energy Center
Sep 28 Fri	Pittsburgh Pirates vs. Cincinnati Reds PNC Park
Sep 28 Fri	Mr. Greengenes Altar Bar

Points of interest




Search Results

This is what's new

- Map
- General info
- Upcoming Events
- Points of interest

*The type of information that appears in this panel depends on what you are searching for

Pittsburgh



Pittsburgh is the second-largest city in the U.S. Commonwealth of Pennsylvania and the county seat of Allegheny County. Regionally, it anchors the largest urban area of both Appalachia and the Ohio River Valley. Wikipedia

Founded: November 25, 1758

Area: 58.3 sq miles (151 km²)

Weather: 61°F (16°C), Wind S at 10 mph (16 km/h), 70% Humidity


Local time: 11:26pm Tuesday (EDT)

Population: 307,484 (2011)

Upcoming events

Sep 26 Wed	Pittsburgh Penguins vs. Detroit Red Wings CONSOL Energy Center
Sep 28 Fri	Pittsburgh Pirates vs. Cincinnati Reds PNC Park
Sep 28 Fri	Mr. Greengenes Altar Bar

Points of interest



Search Engines: State of the Art

- **Input:** Strings (typed or via audio), images, ...
- **Public services:**
 - Links to web pages plus mini synopses via GUI
 - Presentations of structured information via GUI excerpts from the Knowledge Vault (previously known as Knowledge Graph)
- **NSA services: ?**
- **Methods:** Information retrieval, machine learning
- **Data:** Grabbed from many resources (win-win suggested):
 - Web, Wikipedia (DBpedia, Wikidata, ...), DBLP, Freebase, ...

Search Engines

- Find documents: Papers, articles, presentations, ...
 - Extremely cool
 - But...
- Hardly any support for interpreting *documents* w.r.t. certain goals (Knowledge Vault is just a start)
- No support for interpreting *data*
- Claim: Standard search engines provide services but copy documents (and possibly data)
- Why can't individuals provide similar services on their document collections and data?

Personalized Information Engines

- Keep data, provide information
- Invite „agents“ to „view“ (i.e., interpret) local documents and data, without giving away all data
- Let agents take away „their“ interpretation of local documents and data (just like in a reference library).
- Doc/data provider benefits from other agents by (automatically) interacting with them
 - Agents should be provided with incentives to have them „share“ their interpretations
- **No GUI-based interaction, but ...
... semantic interaction via agents**

- **Web and Data Science**
 - Module: **Web-Mining Agents**
 - Machine Learning / Data Mining (Wednesdays)
 - Agents / Information Retrieval (Thursdays)
 - Requirements:
 - Algorithms and Data Structures, Logics, Databases, Linear Algebra and Discrete Structures, Stochastics
 - Module: **Foundations of Ontologies and Databases**
- **Web-based Information Systems**
- **Data Management**
 - Mobile and Distributed Databases
 - Semantic Web

Complementary Courses@UzL

- Algorithmics, Logics, and Complexity
- Signal Processing / Computer Vision
- Machine Learning
- Pattern Recognition
- Artificial Neural Networks (Deep Learning)

Web-Mining Agents

Data Mining

Prof. Dr. Ralf Möller
Universität zu Lübeck
Institut für Informationssysteme

Karsten Martiny (Übungen)



Literature

- Stuart Russell, Peter Norvig, Artificial Intelligence – A Modern Approach, Pearson, 2009 (or 2003 ed.)
- Ian H. Witten, Eibe Frank, Mark A. Hall, Data Mining: Practical Machine Learning Tools and Techniques, Morgan Kaufmann, 2011
- Ethem Alpaydin, Introduction to Machine Learning, MIT Press, 2009
- Numerous additional books, presentations, and videos

Why “Learn” ?

- Machine learning is programming computers to optimize a *performance criterion* using example data or “past experience”
- Simple form of data interpretation
- There is no need to “learn” to calculate payrolls
- Learning is used when:
 - Human expertise does not exist (navigating on planet X),
 - Humans are unable to explain their expertise (speech recognition)
 - Solution changes in time (routing on a computer network)
 - Solution needs to be adapted to particular cases (user biometrics)

What We Talk About When We Talk About “Learning”

- Learning general models from data of particular examples
- Data might be cheap and abundant:
Data warehouse (data mart) maintained by company
- Example in retail: Customer transactions to consumer behavior:
People who bought “Da Vinci Code” also bought “The Five People You Meet in Heaven” (www.amazon.com)
- Build a model that is *a good and useful approximation* to the data

Data Mining

Application of machine learning methods to large databases is called “Data mining”.

- Retail: Market basket analysis, customer relationship management (CRM, also relevant for wholesale)
- Finance: Credit scoring, fraud detection
- Manufacturing: Optimization, troubleshooting
- Medicine: Medical diagnosis
- Telecommunications: Quality of service optimization
- Bioinformatics: Sequence or structural motifs, alignment
- Web mining: Search engines
- ...

What is Machine Learning?

- Optimize a performance criterion using example data or past experience.
- Role of Statistics: Building mathematical models, core task is inference from a sample
- Role of Computer Science: Efficient algorithms to
 - Solve the optimization problem
 - Representing and evaluating the model for inference

Sample of ML Applications

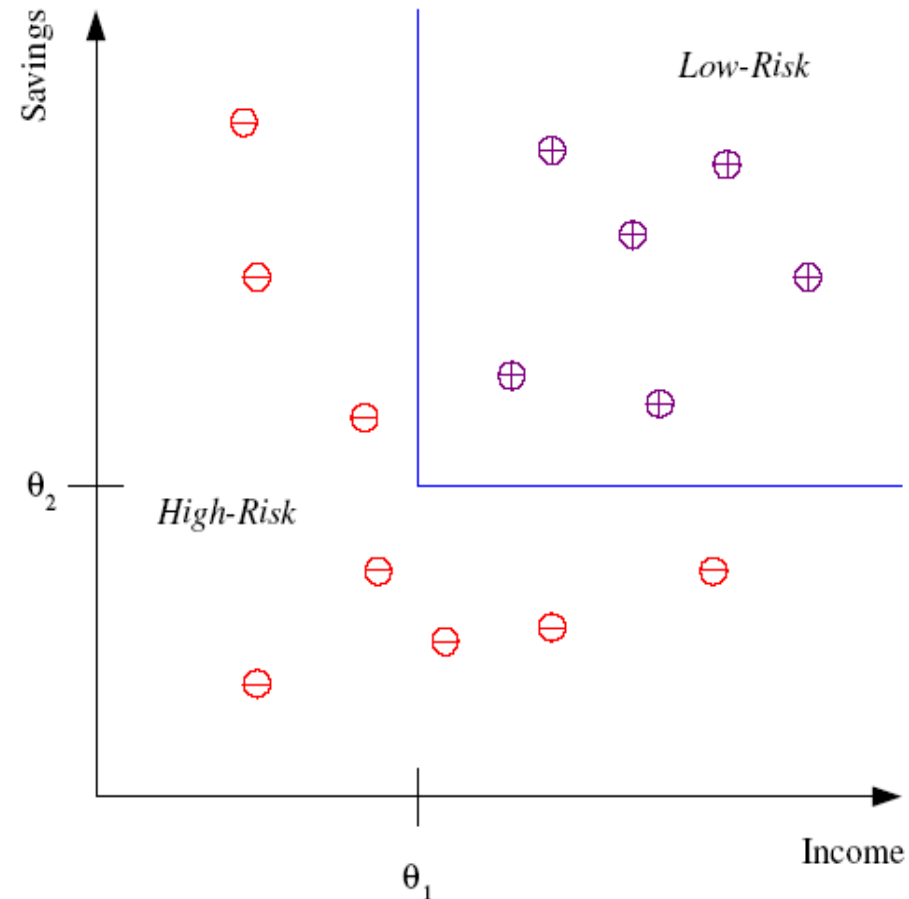
- Learning Associations
- Supervised Learning
 - Classification
 - Regression
- Unsupervised Learning
- Reinforcement Learning

Learning Associations

- Basket analysis:
 $P(Y | X)$ probability that somebody who buys X also buys Y where X and Y are products/services.
Example: $P(\text{chips} | \text{beer}) = 0.7$
- If we know more about customers or make a distinction among them:
 - $P(Y | X, D)$
where D is the customer profile (age, gender, marital status, ...)
 - In case of a web portal, items correspond to links to be shown/prepared/downloaded in advance

Classification

- Example: Credit scoring
- Differentiating between **low-risk** and **high-risk** customers from their *income* and *savings*



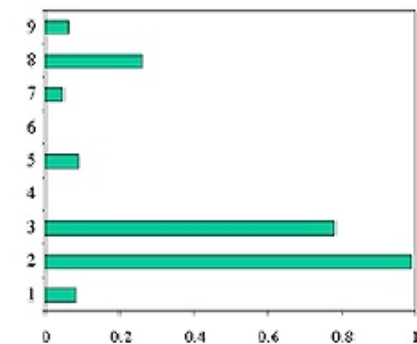
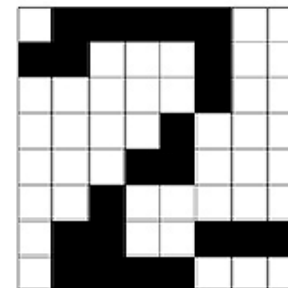
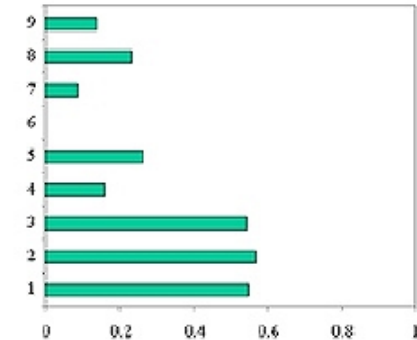
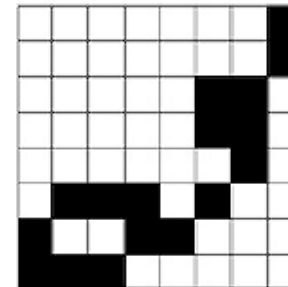
Discriminant: IF $income > \theta_1$ AND $savings > \theta_2$
THEN **low-risk** ELSE **high-risk**

Classification: Applications

- Aka Pattern recognition
- Character recognition: Different handwriting styles.
- Face recognition: Pose, lighting, occlusion (glasses, beard), make-up, hair style
- Speech recognition: Temporal dependency
 - Use of a dictionary for the syntax of the language
 - Sensor fusion: Combine multiple modalities; eg, visual (lip image) and acoustic for speech
- Medical diagnosis: From symptoms to illnesses
- Brainwave understanding: From signals to “states” of thought
- Reading text:
- ...

Character Recognition

Want to learn how to recognize characters, even if written in different ways by different people



Face Recognition

Training examples of a person

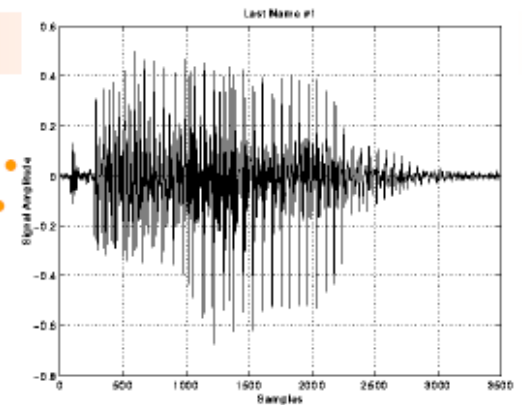


Test images



AT&T Laboratories, Cambridge UK

Example Pattern Recognition: Speech Recognition



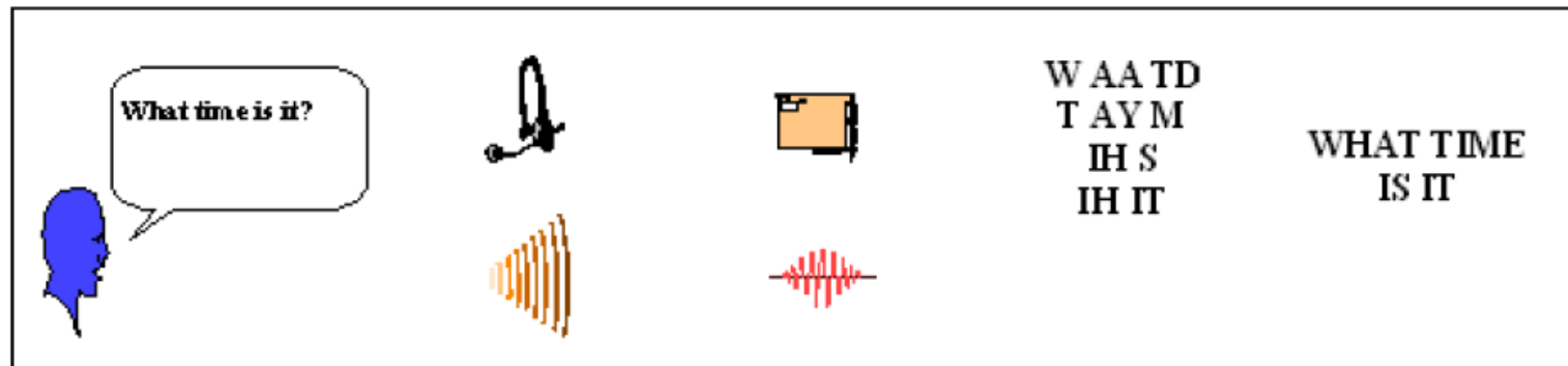
USER

MICROPHONE

SOUND CARD

**SPEECH
RECOGNITION
ENGINE**

**SPEECH-AWARE
APPLICATION**



User speaks into
the microphone.

Microphone captures
sound waves and
generates electrical
impulses.

Sound card
converts
acoustical signal
to digital signal.

Speech recognition
engine converts
digital signal to
phonemes, then
words.

Application
processes
words as text
input.

Medical diagnosis

Inputs: relevant info
about patient, symptoms,
test results, etc.

Output: Expected illness
or risk factors

MEDCAL Risk CHD

PAUL TYERMAN

Age: 61

[Risk Score](#) [Date recorded](#) 18/08/2000

[Cancel](#)

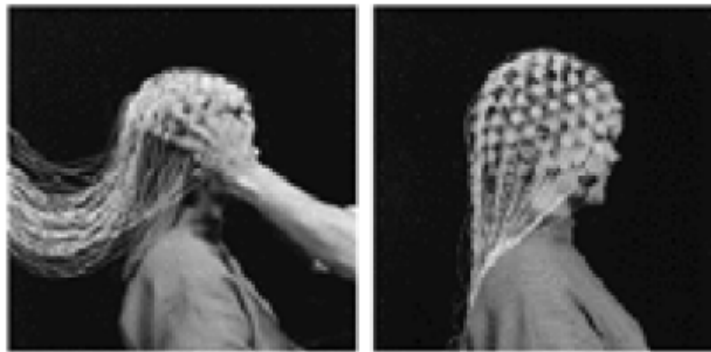
[Finish](#)

24	Blood Pressure	180/95	Exercise Advice	18/08/2000
22	Body Mass Index	29.3	Diet Advice	18/08/2000
20	Smoking	20+	Smoke Advice	18/08/2000
18	Alcohol	13	Drink Advice	No
16	Salt	Not Added	Not printed	Not printed
14	Cholesterol	6.0	Not printed	Not printed
12	HDL/Total ratio	17%	Assessor Number	2
10	Triglycerides	2.0		
8	Diabetic	No		
6	Diabetic relative	Yes		
4	Enlarged heart	No		
2	MI or Angina	No		
0	Family history	40-49		

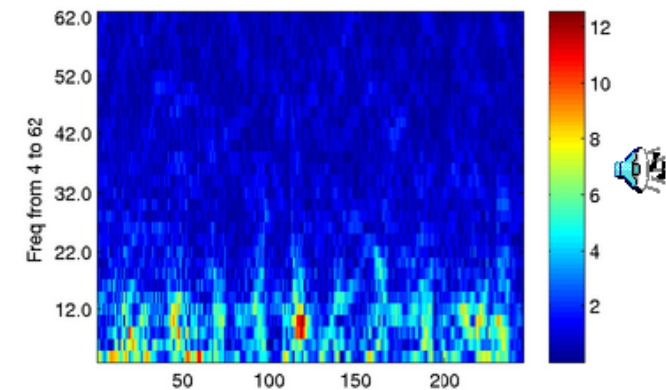
54%

Example Pattern Recognition: Interpreting Brainwaves

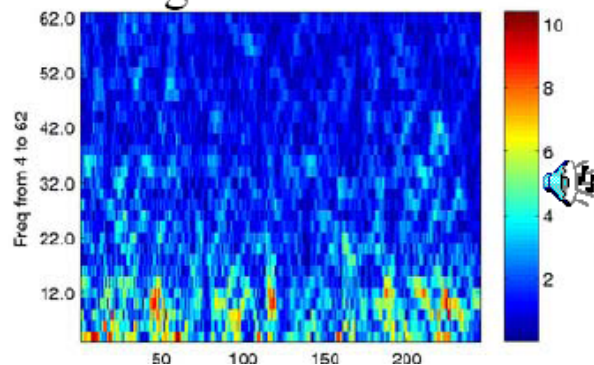
EEG electrodes reading brain waves:



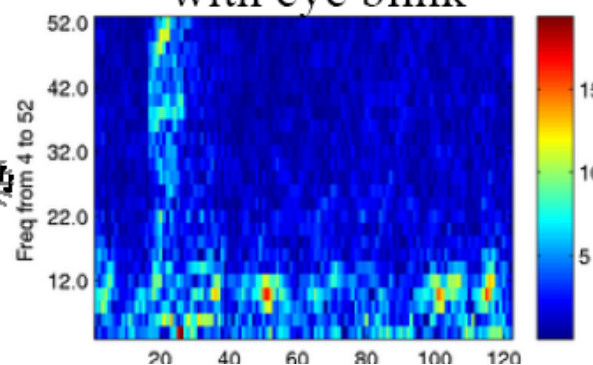
■ Rotation task, left brain



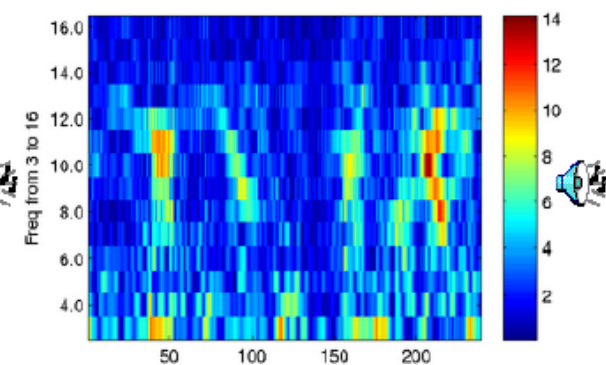
■ Rotation task, right brain



■ Resting task, with eye blink



■ Counting task



Example Pattern Recognition:

Reading text

■ Can you read this?

□ Aircndcog to a rseerhcaer at Cbiardmge Urensvitiy, it dsoen't mtetar in waht oderr the letrtes in a wrod are, the olny ipnaotmrt tihng is taht the fsrit and lsat lteter be at the rgiht plcae. The rset can be a toatl mses and you can slitl raed it wutohit porlebm. Tehy spectluae taht tihs is bseuace the hmaun mnid deos not raed erevy leettr by iesltf but the wrod as a whloe. Wtehehr tihs is ture or not is a ponit of deabte.

■ Clearly, the brain has learned syntax and semantics of language, including contextual dependencies, to make sense of of this 😊

■ For fun: Here's a web page where you can create your own jumbled text: <http://www.stevesachs.com/jumbler.cgi>

Regression

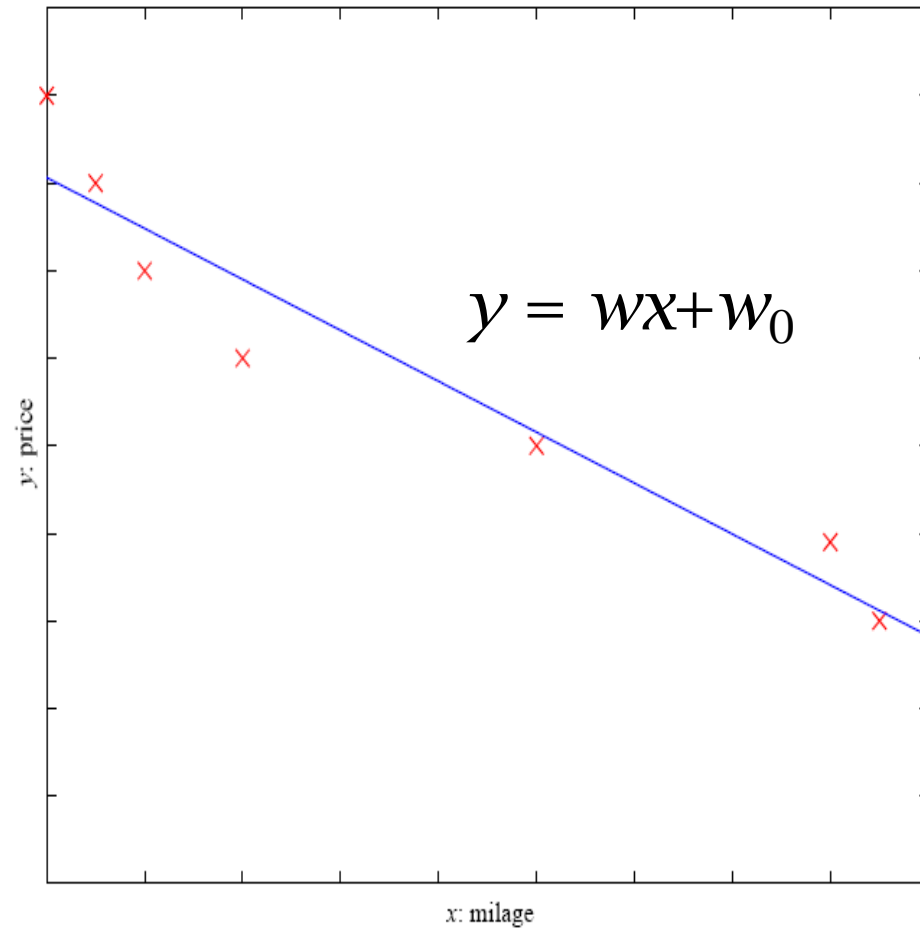
- Example: Price of a used plane
- x : plane attribute

y : price

$$y = g(x | \theta)$$

$g()$ model,

θ parameters



Supervised Learning: Uses

- Prediction of future cases: Use the rule to predict the output for future inputs
- Knowledge extraction: The rule is easy to understand
- Compression: The rule is simpler than the data it explains
- Outlier detection: Exceptions that are not covered by the rule, e.g., fraud

Unsupervised Learning

- Learning “what normally happens”
- No output (we do not know the right answer)
- Clustering: Grouping similar instances
- Example applications
 - Customer segmentation in CRM
 - Company may have different marketing approaches for different groupings of customers
 - Image compression: Color quantization
 - Instead of using 24 bits to represent 16 million colors, reduce to 6 bits and 64 colors, if the image only uses those 64 colors
 - Bioinformatics: Learning motifs (sequences of amino acids in proteins)
 - Document classification in unknown domains

Reinforcement Learning

- Learning a policy: A sequence of actions/outputs
- No supervised output but delayed reward
- Credit assignment problem
- Game playing
- Robot in a maze
- Multiple agents, partial observability, ...

An Extended Example

- “Sorting incoming Fish on a conveyor according to species using optical sensing”

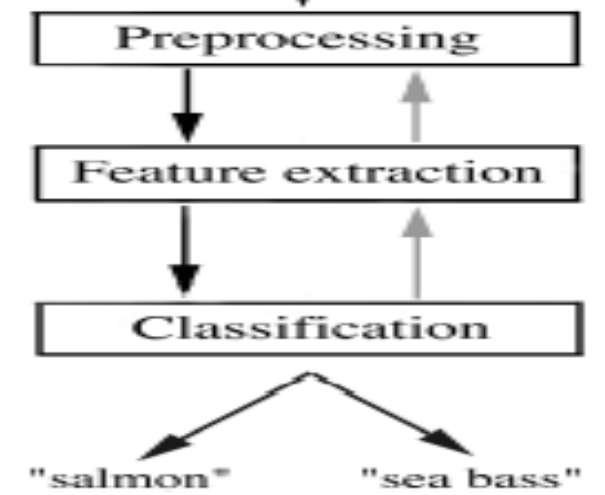


Problem Analysis

- Set up a camera and take some sample images to extract features
 - Length
 - Lightness
 - Width
 - Number and shape of fins
 - Position of the mouth, etc...
- This is the set of all suggested features to explore for use in our classifier!

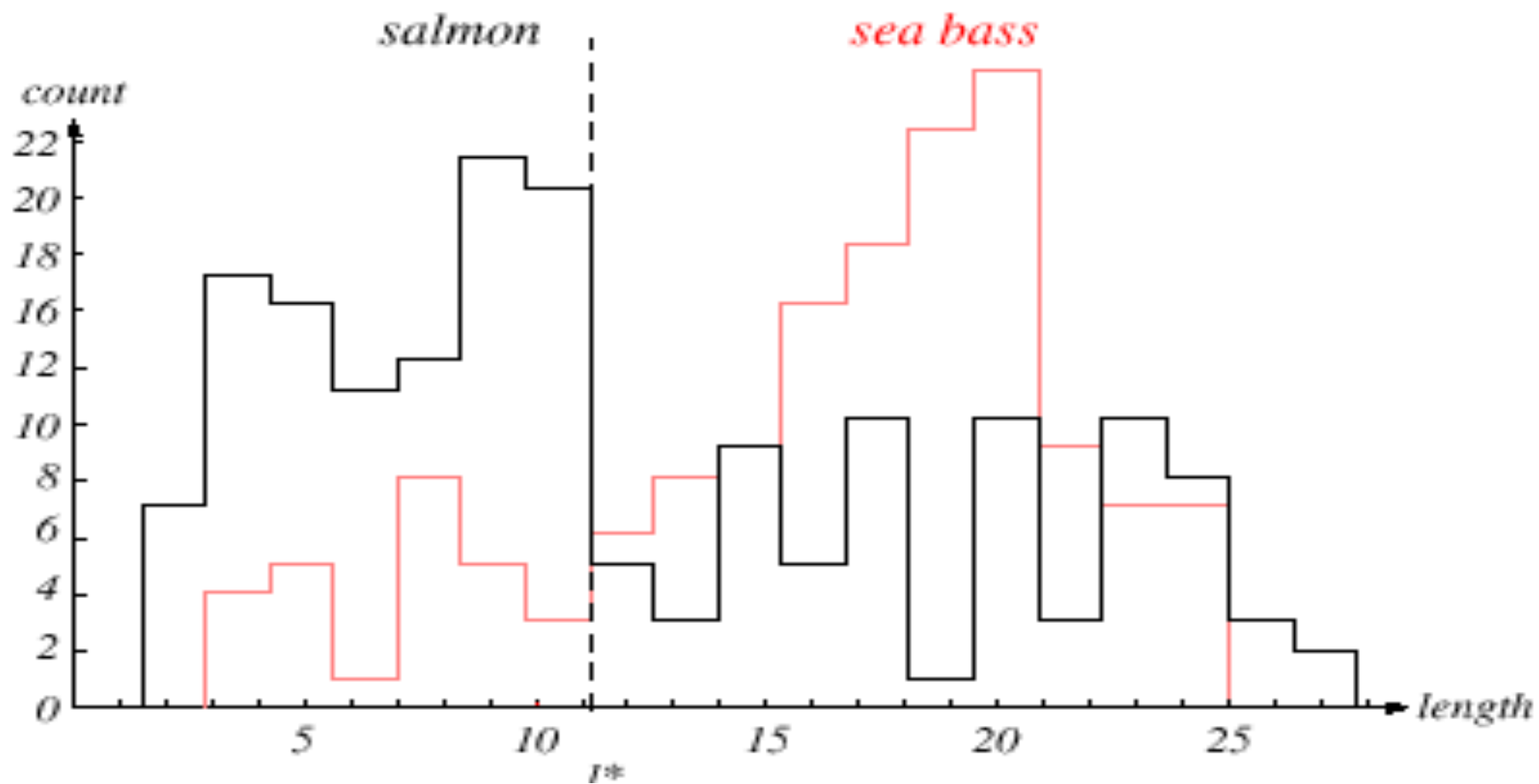
Preprocessing

- Use a segmentation operation to isolate fishes from one another and from the background
- Information from a single fish is sent to a feature extractor whose purpose is to reduce the data by measuring certain features
- The features are passed to a classifier



Classification

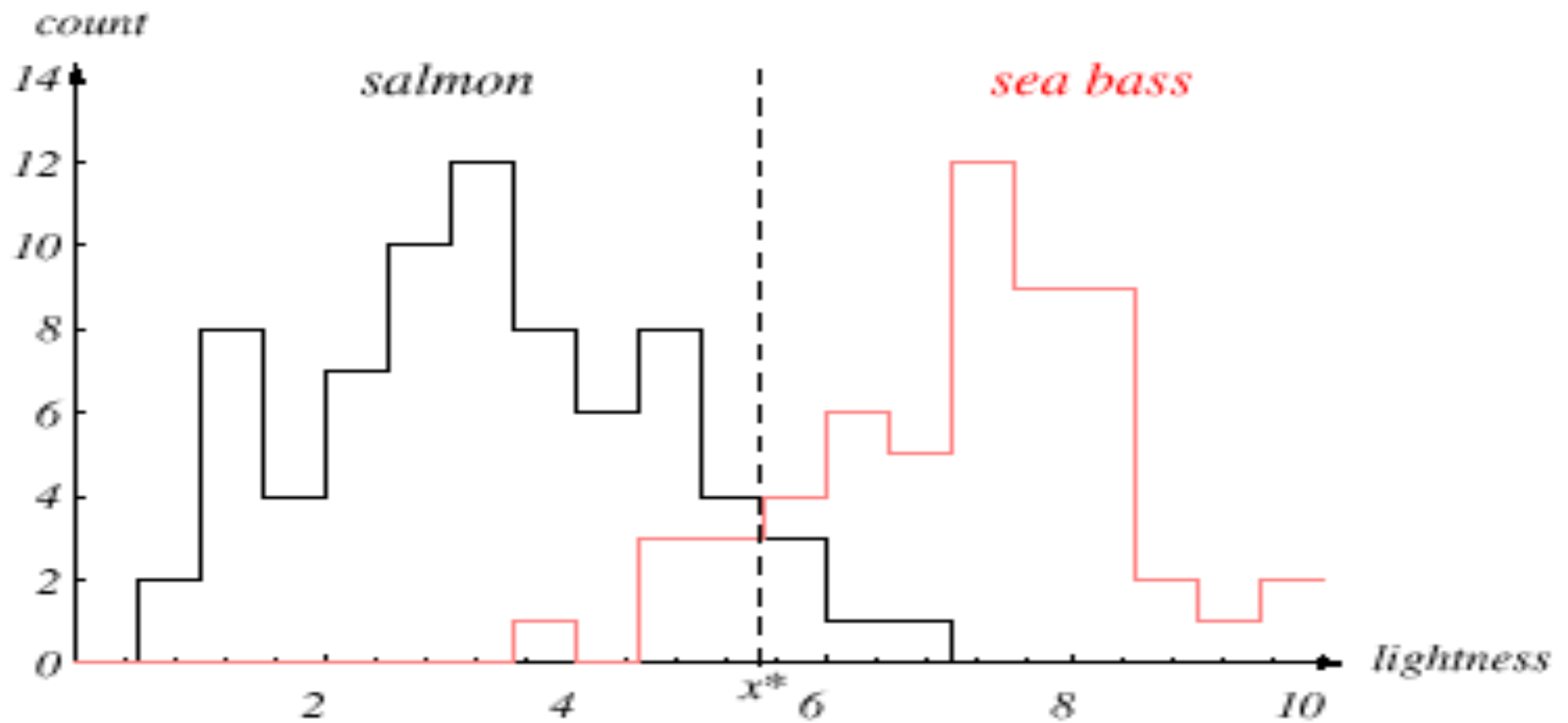
- Now we need (expert) information to find features that enables us to distinguish the species.
- “Select the length of the fish as a possible feature for discrimination”



The **length** is a poor feature alone!

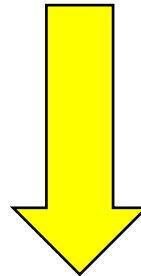
→ Cost of decision

Select the **lightness** as a possible feature.



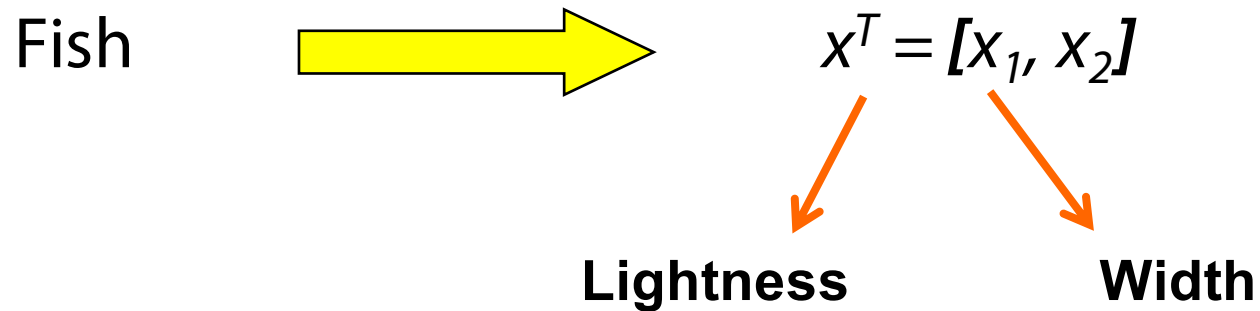
Threshold decision boundary and cost relationship

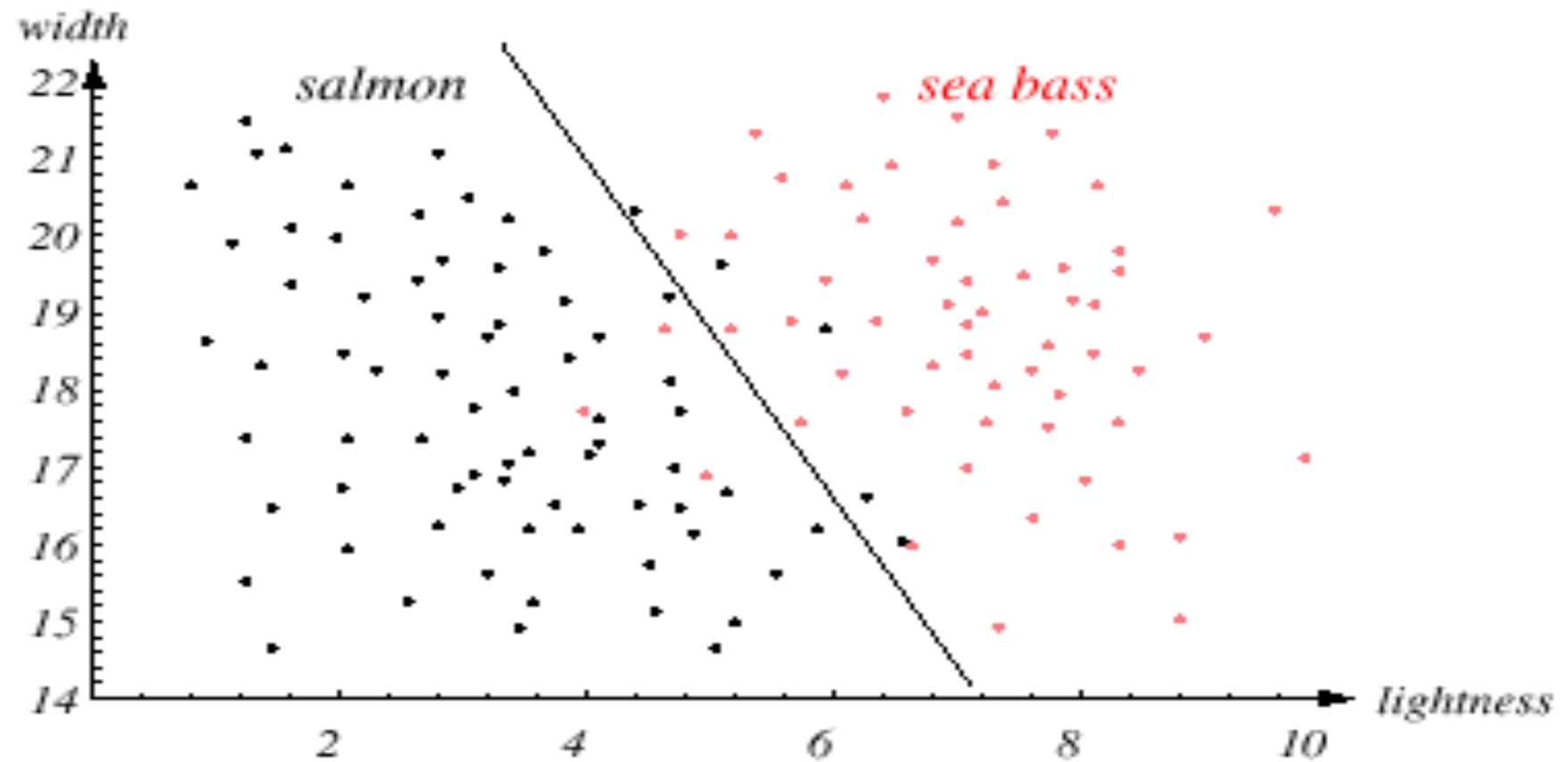
- Move our decision boundary toward smaller values of lightness in order to minimize the cost (reduce the number of sea bass that are classified salmon!)



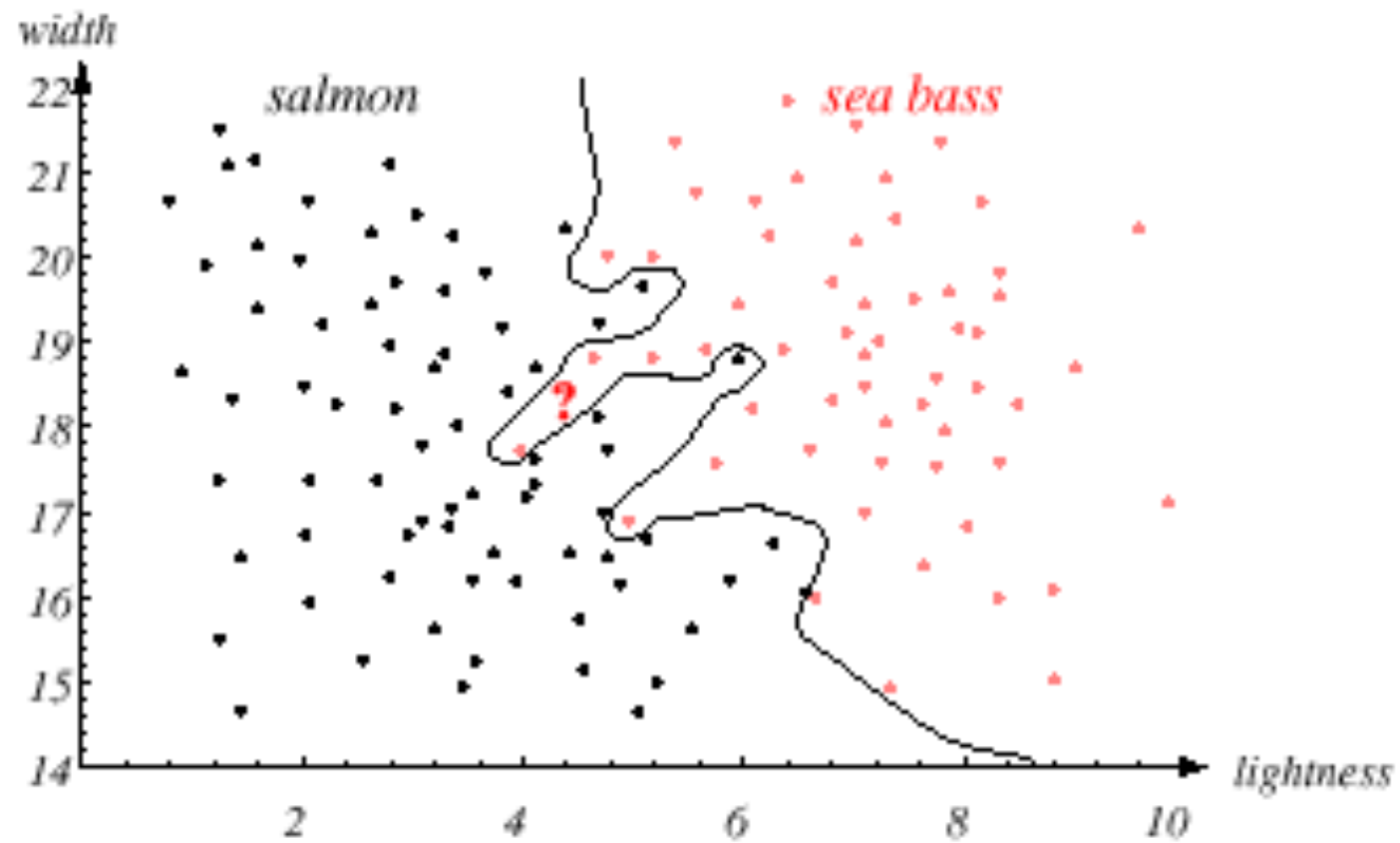
Task of decision theory

Adopt the lightness and add the width of the fish

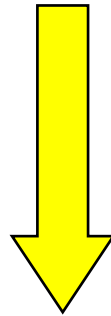




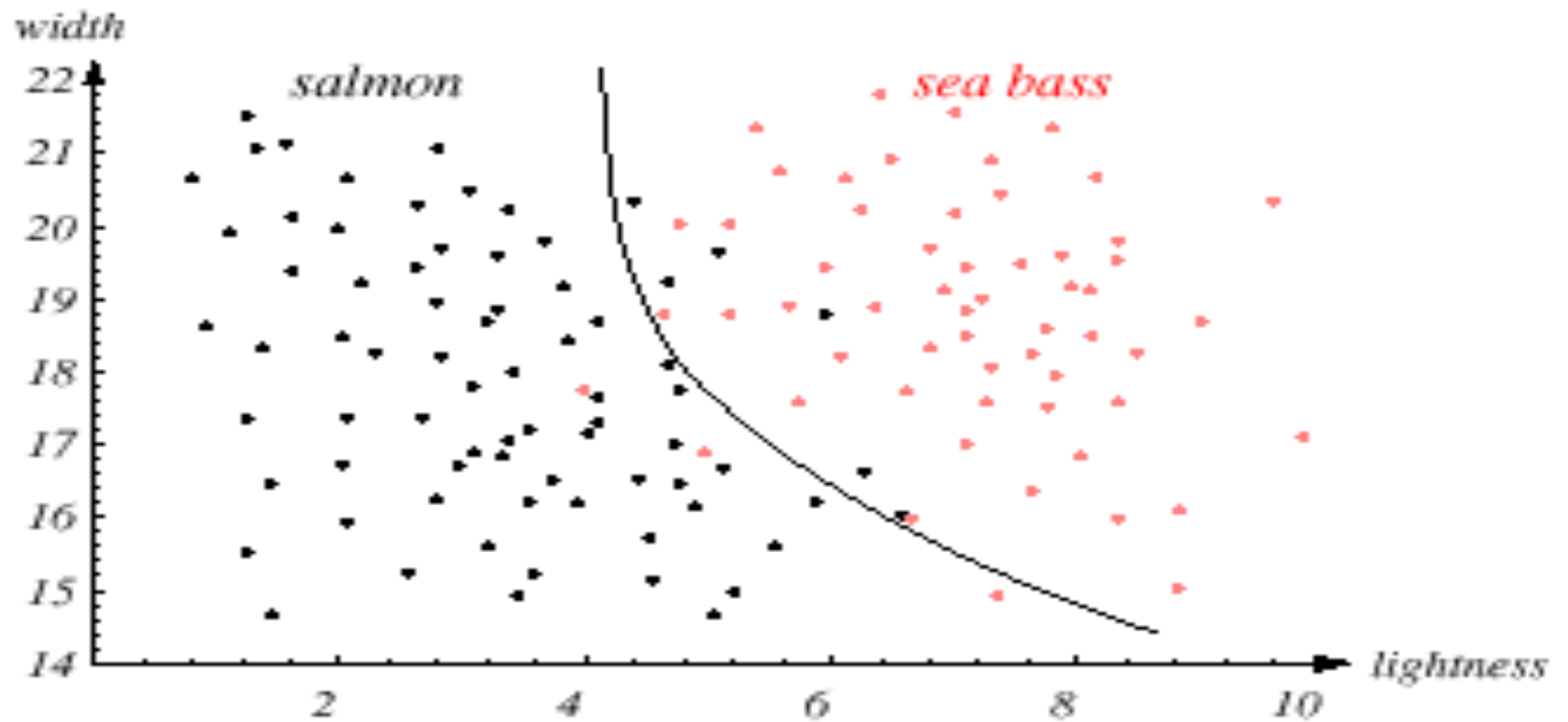
-
- We might add other features that are not correlated with the ones we already have.
 - Precaution should be taken not to reduce the performance by adding such “noisy features”
 - Ideally, the best decision boundary should be the one which provides an optimal performance



However, our satisfaction is premature because the central aim of designing a classifier is to correctly classify novel input

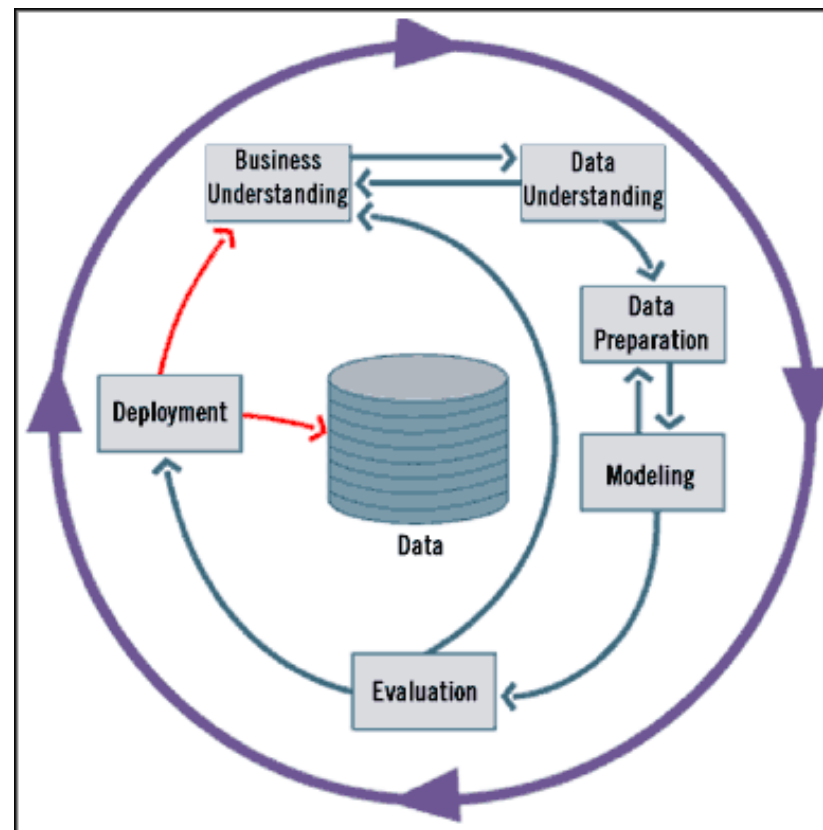


Issue of generalization!



Standard data mining life cycle

- It is an iterative process with phase dependencies
- Consists of six (6) phases:



Phases (1)

- Business Understanding
 - Understand project objectives and requirements
 - Formulation of a data mining problem definition
- Data Understanding
 - Data collection
 - Evaluate the quality of the data
 - Perform exploratory data analysis
- Data Preparation
 - Clean, prepare, integrate, and transform the data
 - **Select** appropriate attributes and variables

Phases (2)

- Modeling
 - Select and apply appropriate modeling techniques
 - Calibrate/learn model parameters to optimize results
 - If necessary, return to data preparation phase to satisfy model's data format
- Evaluation
 - Determine if model satisfies objectives set in phase 1
 - Identify business issues that have not been addressed
- Deployment
 - Organize and present the model to the “user”
 - Put model into practice
 - Set up for continuous mining of the data

Fallacies of Data Mining (1)

- Fallacy 1: There are data mining tools that automatically find the answers to our problem
 - Reality: There are no automatic tools that will solve your problems “while you wait”
- Fallacy 2: The DM process requires little human intervention
 - Reality: The DM process require human intervention in all its phases, including updating and evaluating the model by human experts
- Fallacy 3: Data mining have a quick ROI
 - Reality: It depends on the startup costs, personnel costs, data source costs, and so on

Fallacies of Data Mining (2)

- Fallacy 4: DM tools are easy to use
 - Reality: Analysts must be familiar with the model
- Fallacy 5: DM will identify the causes to the business problem
 - Reality: DM tools only identify patterns in your data, analysts must identify the cause
- Fallacy 6: Data mining will clean up a data repository automatically
 - Reality: Sequence of transformation tasks must be defined by an analysts during early DM phases

Remember

- Problems suitable for Data Mining:
 - Require to discover knowledge to make right decisions
 - Current solutions are not adequate
 - Expected high-payoff for the right decisions
 - Have accessible, sufficient, and relevant data
 - Have a changing environment
- IMPORTANT:
 - **ENSURE privacy if personal data is used!**
 - **Not every data mining application is successful!**

Overview

Supervised Learning



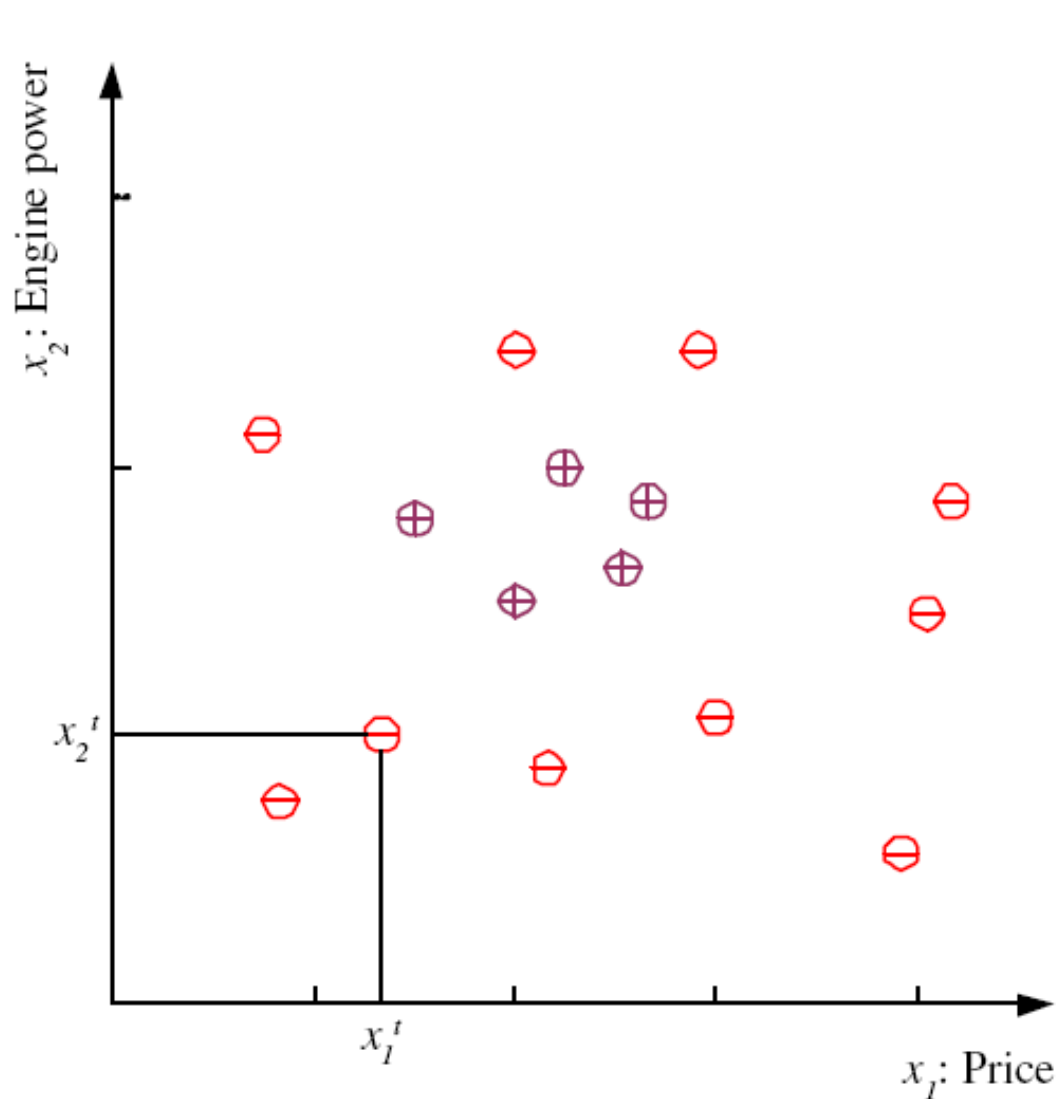
Learning a Class from Examples

- Class C of a “family car”
 - Prediction: Is car x a family car?
 - Knowledge extraction: What do people expect from a family car?
- Output:

Positive (+) and negative (–) examples
- Input representation:

x_1 : price, x_2 : engine power

Training set \mathcal{X}

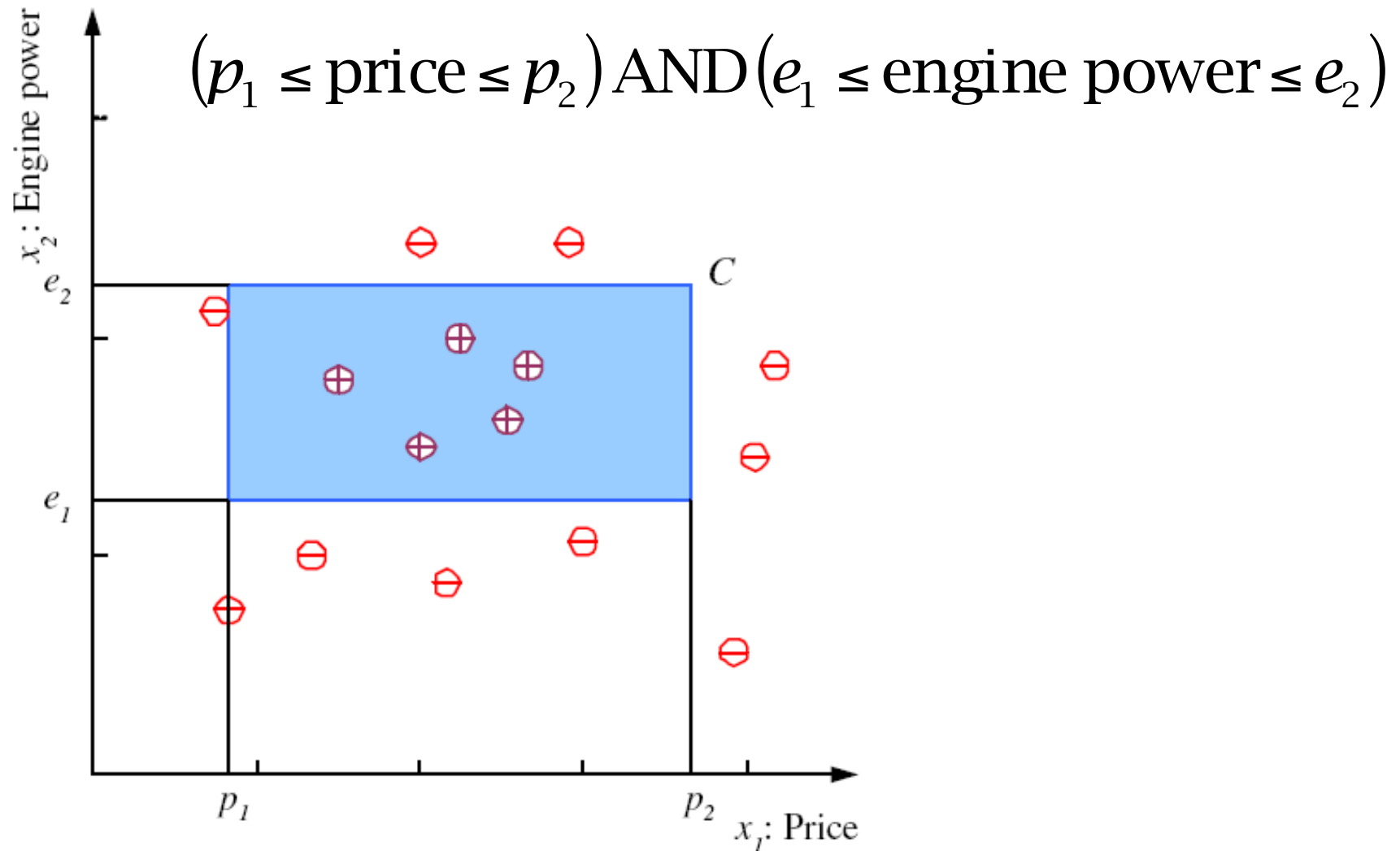


$$\mathcal{X} = \{\mathbf{x}^t, r^t\}_{t=1}^N$$

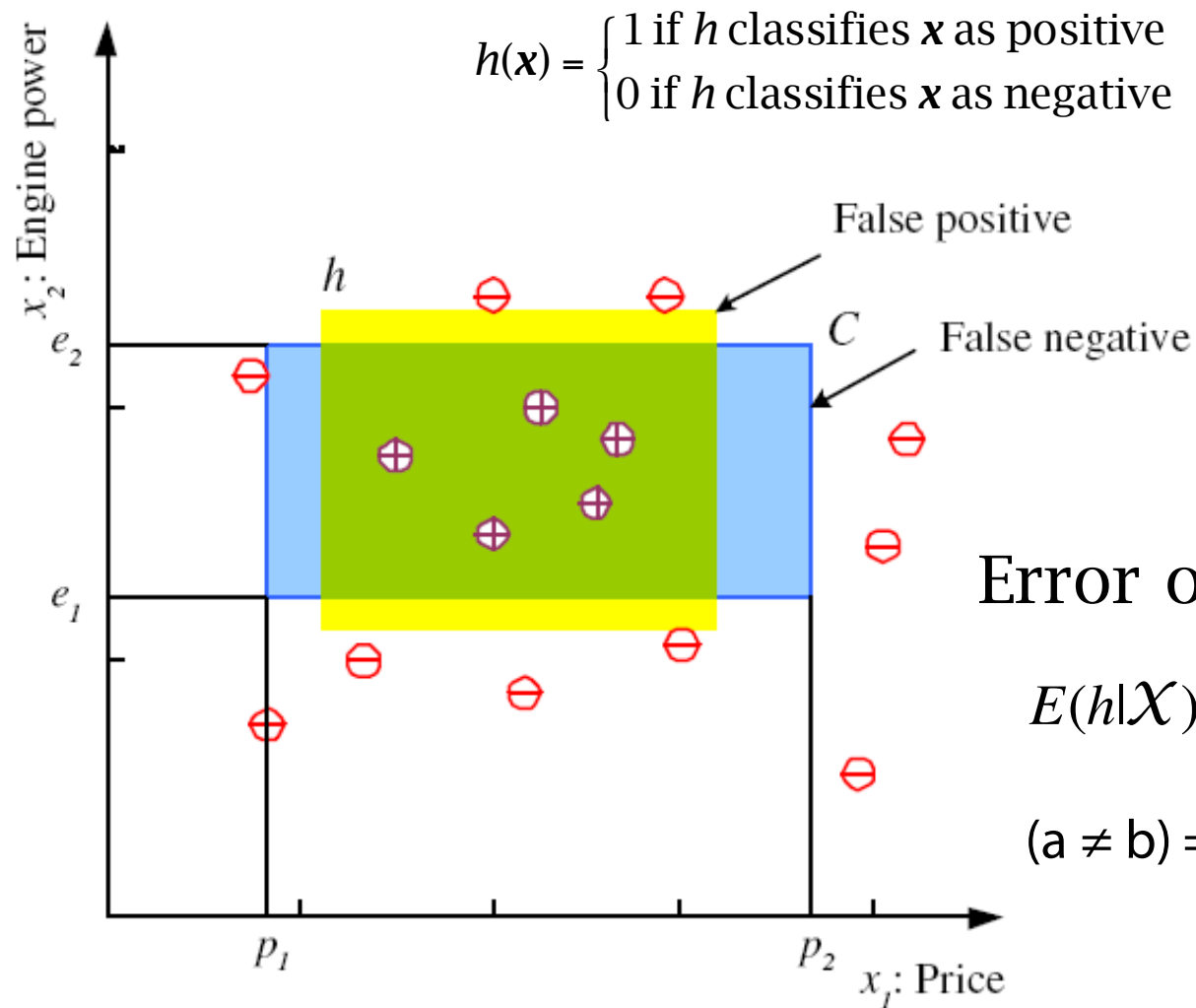
$$r = \begin{cases} 1 & \text{if } \mathbf{x} \text{ is positive} \\ 0 & \text{if } \mathbf{x} \text{ is negative} \end{cases}$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

Class C



Hypothesis class \mathcal{H}

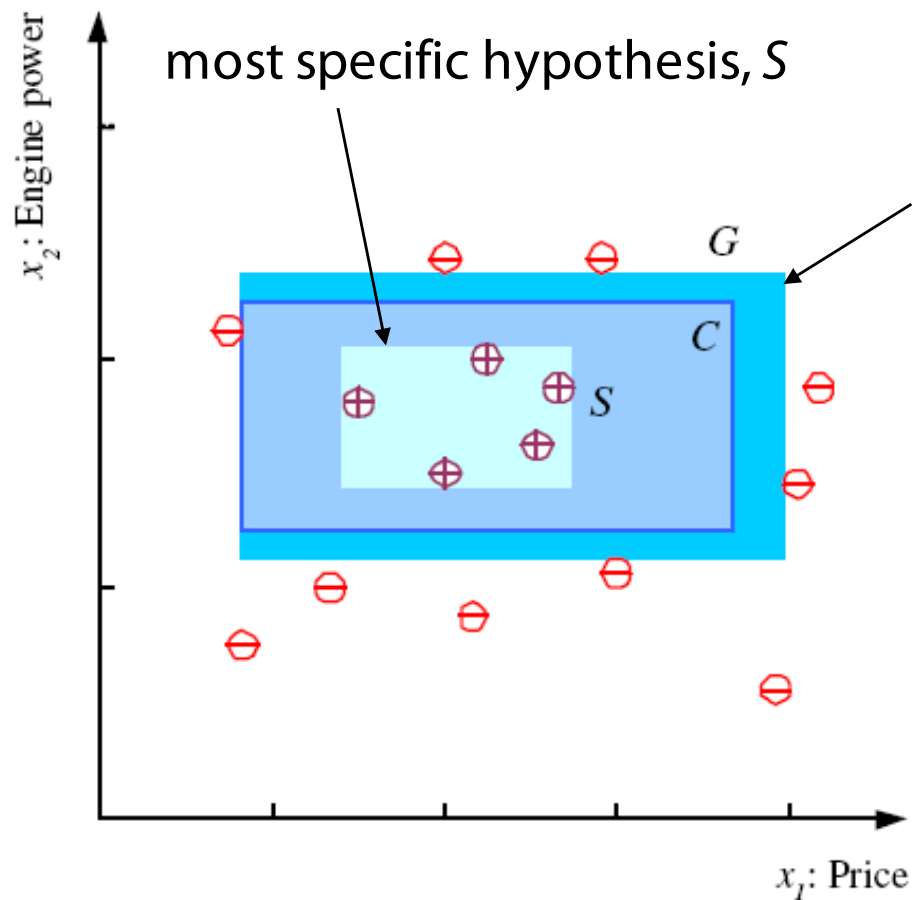


Error of h on \mathcal{H}

$$E(h|\mathcal{X}) = (1/N) \sum_{t=1}^N (h(\mathbf{x}^t) \neq r^t)$$

$(a \neq b) = 1$ if \neq , 0 otherwise

S, G, and the Version Space



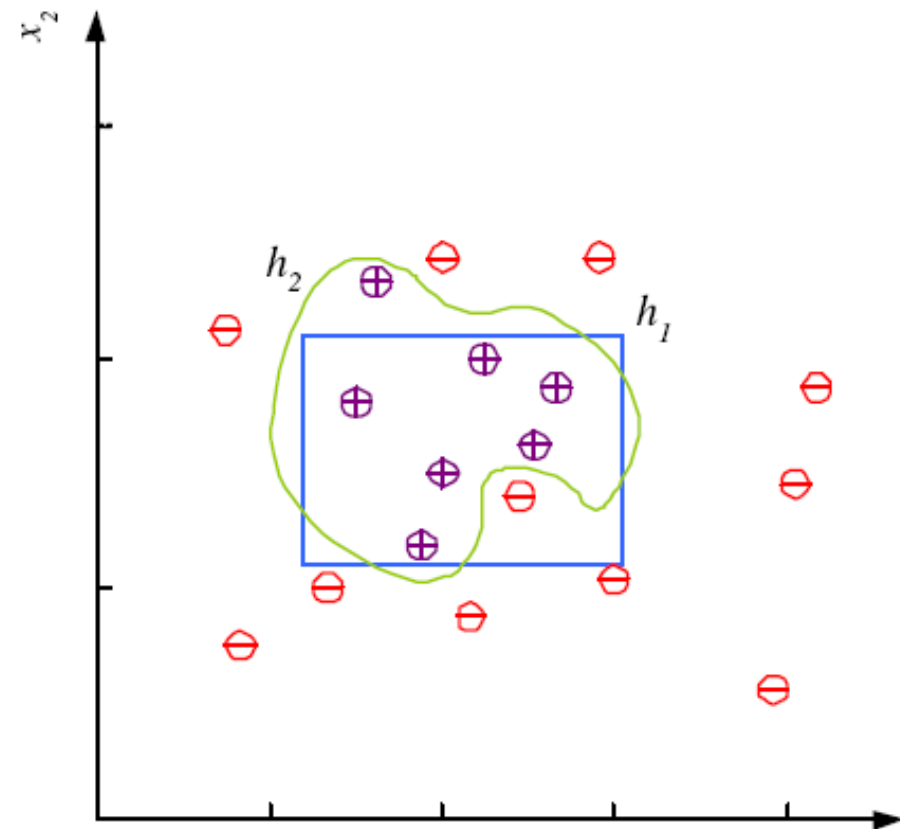
$h \in H$, between S and G is consistent and make up the version space

(Mitchell, 1997)

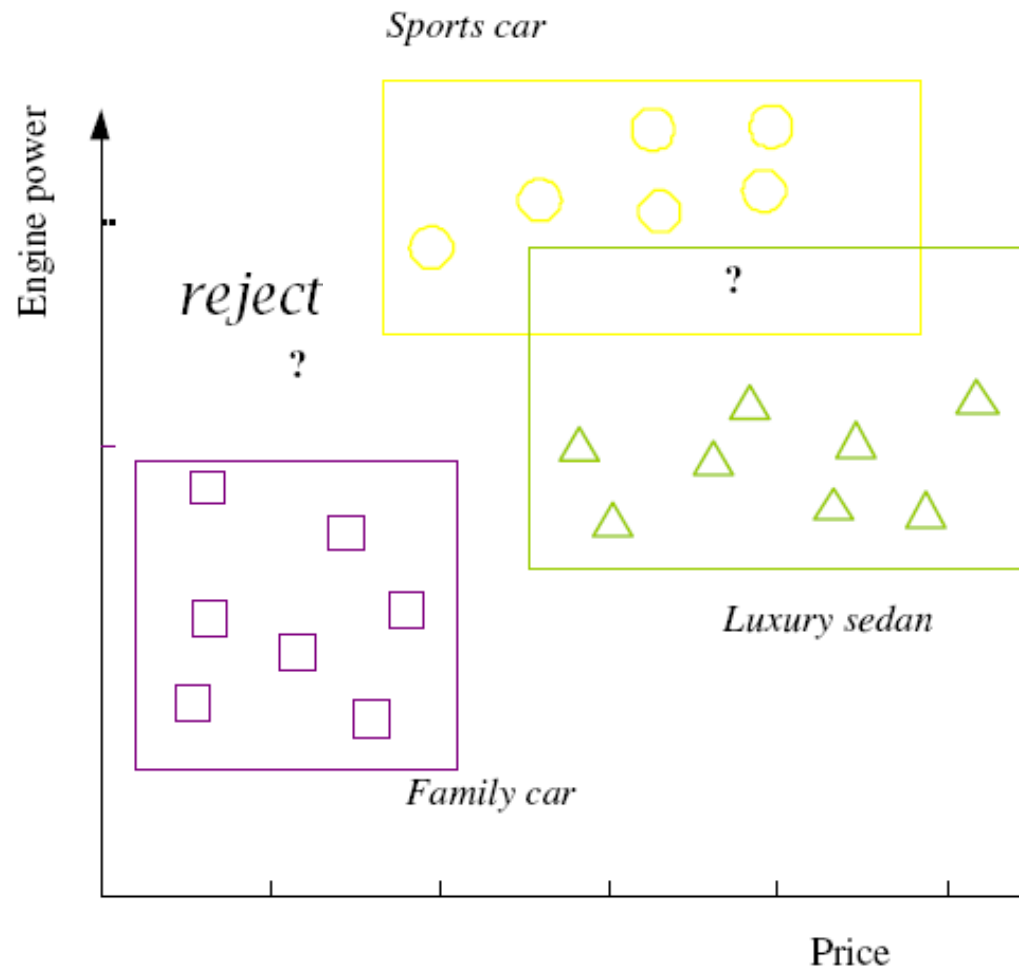
Noise and Model Complexity

Use the simpler one because

- Simpler to use
(lower computational complexity)
- Easier to train (lower space complexity)
- Easier to explain
(more interpretable)
- Generalizes better (lower variance - Occam's razor)



Multiple Classes, C_i $i=1,...,K$



$$\mathcal{X} = \{\mathbf{x}^t, r^t\}_{t=1}^N$$

$$r_i^t = \begin{cases} 1 & \text{if } \mathbf{x}^t \in C_i \\ 0 & \text{if } \mathbf{x}^t \in C_j, j \neq i \end{cases}$$

Train hypotheses

$h_i(\mathbf{x}), i = 1, \dots, K:$

$$h_i(\mathbf{x}^t) = \begin{cases} 1 & \text{if } \mathbf{x}^t \in C_i \\ 0 & \text{if } \mathbf{x}^t \in C_j, j \neq i \end{cases}$$

Regression

$$\mathcal{X} = \{x^t, r^t\}_{t=1}^N$$

$$r^t \in \mathfrak{R}$$

$$r^t = f(x^t)$$

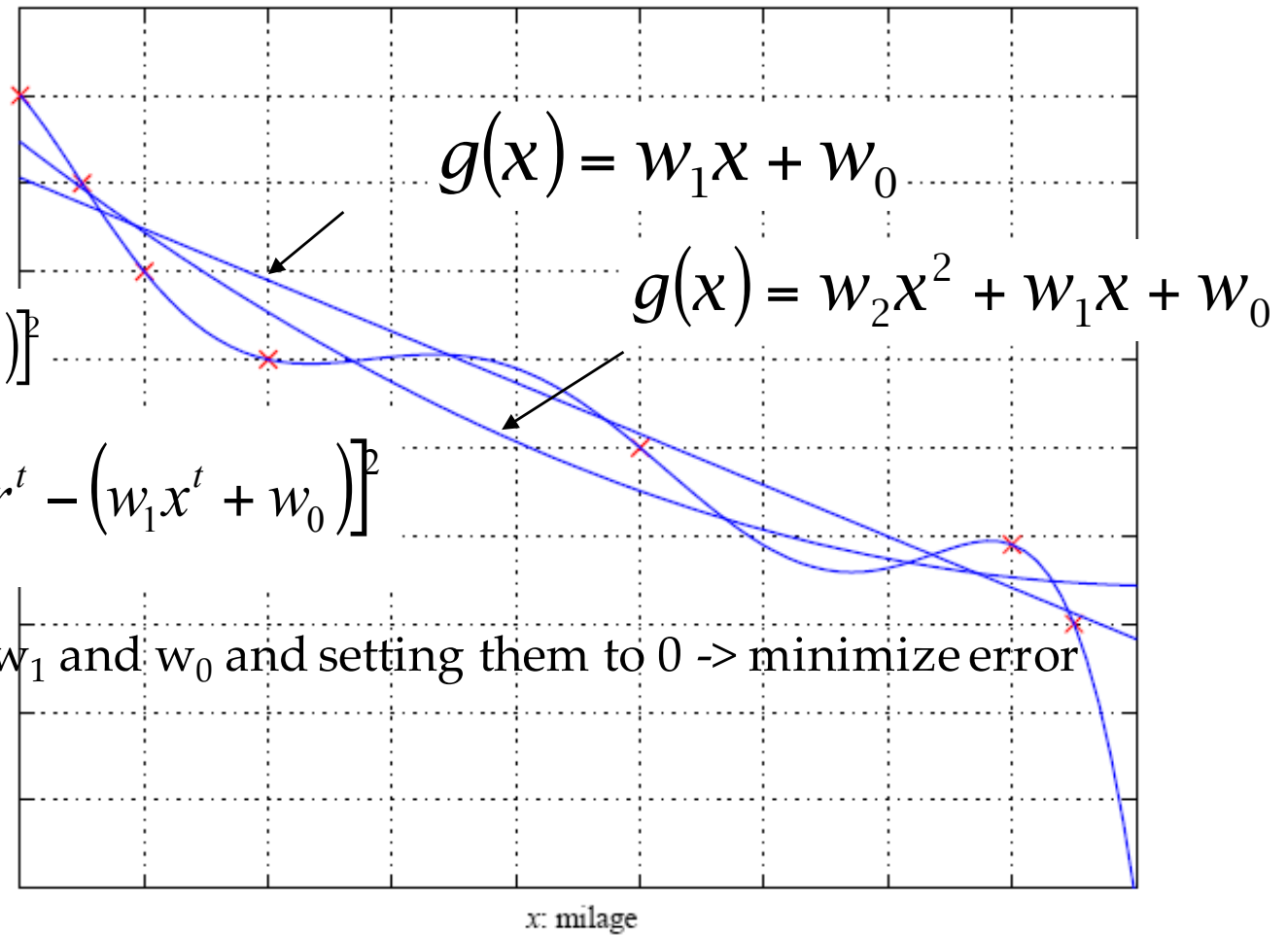
$$E(g | \mathcal{X}) = \frac{1}{N} \sum_{t=1}^N [r^t - g(x^t)]^2$$

$$E(w_1, w_0 | \mathcal{X}) = \frac{1}{N} \sum_{t=1}^N [r^t - (w_1 x^t + w_0)]^2$$

Partial derivatives of E w.r.t w_1 and w_0 and setting them to 0 -> minimize error

$$w_1 = \frac{\sum_t x^t r^t - \bar{x} \bar{r} N}{\sum_t (x^t)^2 - N \bar{x}^2}$$

$$w_0 = \bar{r} - w_1 \bar{x}$$



Model Selection & Generalization

- Learning is an **ill-posed problem**;
data is not sufficient to find a unique solution
- The need for **inductive bias**, assumptions about **H**
- **Generalization**: How well a model performs on new data
- Overfitting: **H** more complex than **C** or **f**
- Underfitting: **H** less complex than **C** or **f**

Triple Trade-Off

There is a trade-off between three factors (Dietterich, 2003):

1. Complexity of \mathcal{H} , $c(\mathcal{H})$,
 2. Training set size, N ,
 3. Generalization error, E , on new data
- As N , E^-
 - As $c(\mathcal{H})$, first $E \downarrow$ and then E

Cross-Validation

- To estimate generalization error, we need data unseen during training. We split the data as
 - Training set (50%)
 - Validation set (25%)
 - Test (publication) set (25%)
- Resampling when there is few data

Dimensions of a Supervised Learner

1. Model: $g(\mathbf{x} \mid \theta)$
2. Loss function: $E(\theta \mid \mathcal{X}) = \sum_t L(r^t, g(\mathbf{x}^t \mid \theta))$
3. Optimization procedure: $\theta^* = \arg \min_{\theta} E(\theta \mid \mathcal{X})$

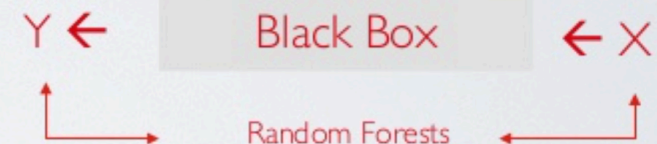
Data Models vs. Algorithmic Models

Data Modeling

VS.

Algorithmic Modeling

$$Y \leftarrow F(X, \text{random noise, parameters})$$



We understand the world

How well 'my data model' works
Statisticians, Data Analysts, Data Miners
Linear Regression
Logistic Regression
Known Distributions
Confidence Intervals
Predictor Variables & Goodness of Fit

We don't understand the world

The world produces data in a black-box
Data Scientists
Machine Learning, AI & Neural Nets
Random Forests, SVM, GBT
Unknown Multivariate Distributions
Iterative
Predictive Accuracy

"Statistical Modeling: The Two Cultures" Leo Breiman, 2001