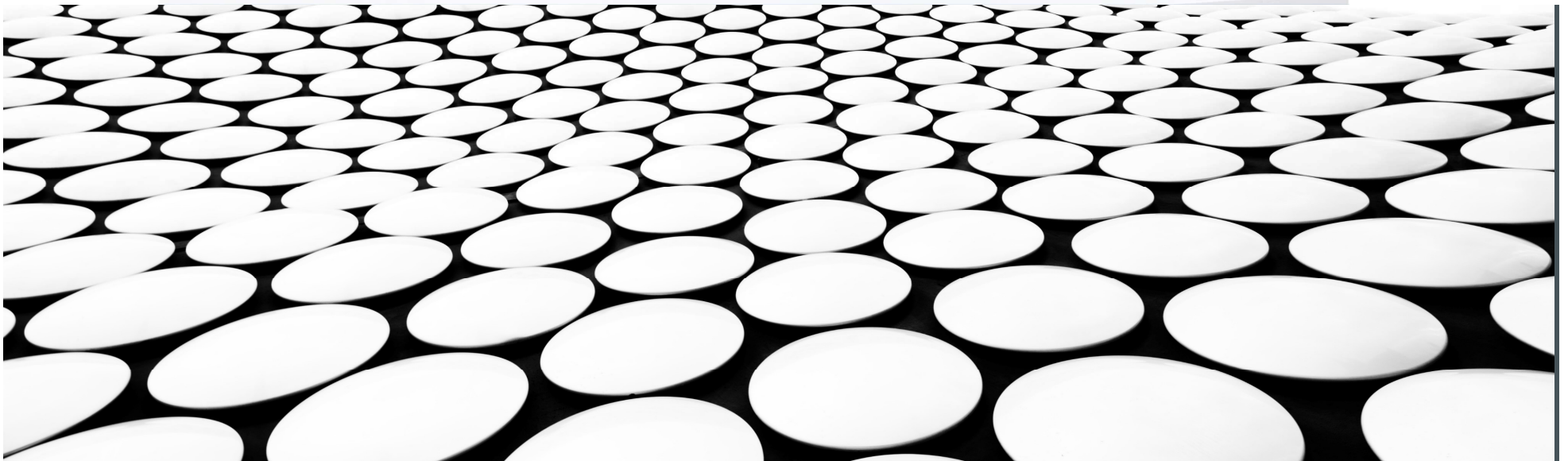

INTRODUCTION TO STATISTICAL LEARNING THEORY

ROBERTO BASILI, UNIVERSITÀ DI ROMA, TOR VERGATA

WEB MINING & RETRIEVAL, MARCH 2023



A COURSE ROADMAP

Linguistic knowledge in kernel machines

Part I

Transformers and Applications

Part III

Neural Encoding

Part II



OUTLINE

- Language Processing: between knowledge and structures
- Statistical Learning Theory and SVMs
- Kernel Machines
 - Non linearity and kernels
 - Sequence kernels
 - Tree Kernels
- Learning under knowledge constraints
- Embedding Knowledge in NNs
- Laboratory: A USE CASE – Machine Learning for Question Answering



LANGUAGE PROCESSING ... A PROLOGUE

SEMANTICS, OPEN DATA AND NATURAL LANGUAGE

- Web contents, characterized by rich multimedia information, are mostly opaque from a semantic standpoint

Today is 2011年11月13日 星期日 顯示器最佳分辨率1024x768

2011 中國證券金紫荊獎 Golden Bauhinia Awards

首頁 國內 國際 港澳 兩岸 評論 財經 體育 教育 科技 醫藥 娛樂 文化 副刊 軍事 生活 旅遊 圖片 博客

關鍵詞: 欄目: 全部 最近三個月 三個月之前 搜索

滾動新聞:

胡總語特首:防範經濟金融風險

胡錦濤在夏威夷會見出席APEC峰會的曾蔭權。他祝賀香港區議會選舉成功,並充分肯定曾蔭權及港府工作,要求做好經濟金融風險防範。

胡連會登場 共同宣示九二共識

胡錦濤第四度在APEC峰會期間會見連戰。他強調,認同「九二共識」是兩岸開展對話協商的必要前提,也是兩岸關係和平發展的重要基礎。

西藏黨代會高調反「藏獨」

德國作家:外媒錯誤報導西藏
傳媒入日本福島核電站採訪 英國大裁軍 傷兵難倖免
演礦難已30死 13人生還 礦工講述內幕 事故並不意外
范徐麗泰認民望跌最不熱 選委再選60提名表 累積逾千人
聖保羅中學本月底截止招 選委再選60提名表 累積逾千人
民調逆轉 藍高層:國親吵鬧地 秋門訴求多 向藍綠表不滿
世界新七奇觀 亞洲景佔四席 新奇觀選舉黨爭讓
中國實體書店苦苦掙扎求 加入TPP 台密集會談探路
香港人家/蔡仕榮 人生導師 活出自我 香港人家/蔡仕榮
債務危機糾 港ADR幾全線 歐元反彈 兌美元逼近1.38
入世十年/充分對接 華強北最 入世十年/挑戰「二次」
抽身離 工人險生 南亞演說會 暫拘日籍妻

即時新聞

- 組圖/河南全國太極拳錦標賽賽況
- 奧巴馬重申美不支持「台灣獨立」
- 巴基斯坦西北部兩起襲擊 16人死
- 圖文/胡錦濤會見美國總統奧巴馬 (圖)
- 兩岸30對愛侶在廈門集體證婚
- 中日韓衛生部長會議在青島舉行
- 面向中國遊客中英雜誌紐約創刊
- 「CEO聖經」成內地官員考試內容
- 斯特惠:經紀人是勞資談判的障礙
- 香港獲成爲人民幣國際化關鍵角色
- 日學者提出地核物質形態新假說
- 中國影視機備向國際大師「取經」

焦點關注

- 區議會選舉
- 香港特首選舉
- 2011APEC 港果金事件
- 2011施政報告
- 神八天宮对接
- 第七次陳江會
- 李克強訪港
- 9.1衝擊事件
- 中國航母試航
- 辛亥革命百年

http://www.takungpao.com.hk/news/11/11/13/2011_apcec_xgbd-1423309.htm

INFORMATION, WEB AND NATURAL LANGUAGES

Chinese President Hu Jintao (R) shakes hands with Honorary Chairman of the Chinese Kuomintang (KMT) Lien Chan, in Honolulu, Hawaii, the U.S., Nov. 11, 2011.

(Xinhua/Huang Jingwen)

HONOLULU, United States, Nov. 11 (Xinhua) -- Hu Jintao, general secretary of the Central



Chinese President Hu Jintao (R) shakes hands with Honorary Chairman of the Chinese Kuomintang (KMT) Lien Chan, in Honolulu, Hawaii, the U.S., Nov. 11, 2011.
(Xinhua/Huang Jingwen)

HONOLULU, United States, Nov. 11 (Xinhua) -- Hu Jintao, general secretary of the Central



Who is Hu Jintao?

- 3 China in APEC: a mutually beneficial en...
- 4 Night life in Shanghai
- 5 China's 2011 foreign trade to grow 20 p...
- 6 Beijing house prices stumble 5.1 pct as...
- 7 Lama students start school in Tibet Col...
- 8 Police in central China crack phoney ca...
- 9 China-ASEAN cooperation sees notable pr...
- 10 Miao ethnic group celebrates Miao's New...



Hu Jintao



Ricerca

Circa 725.000 risultati (0,09 secondi)

Tutto

Immagini

Mappe

Video

Notizie

Shopping

Più conte

Tutti i ri

Per argomento

Qualsiasi dimensione

Grandi

Medie

icone

Maggiori di...

Dimensioni esatte...

Qualsiasi colore

A colori

Bianco e nero



Qualsiasi tipo

Volti

Foto

Clip art

Disegni

Visual: standard

Mostra dimensioni



CONTENT SEMANTICS AND NATURAL LANGUAGE

- Human languages are the **main carrier of the information involved in knowledge retrieval, communication and exchange** as it is associated to the open Web contents
- **Words and language structures provide all we need** to express concepts, activities, events, abstractions and conceptual relations we usually share through data
- ***“Language is parasitic to knowledge representation languages but the viceversa is not true”*** (Wilks, 2001)
- From Learning to Read to Knowledge Distillation and Management we perform a(n integrated pool of) semantic interpretation task(s) whose automation imply a crucial interest for Data Science.

NATURAL LANGUAGE & AMBIGUITY



Take extra care with children

Dogs must be carried

NATURAL LANGUAGE & AMBIGUITY



- *"Dogs must be carried on this escalator"*

can be interpreted in a number of ways:

- *"All dogs should have a chance to go on this wonderful escalator ride"*
- *"This escalator is for dog-holders only"*
- *"You can't carry your pet on the other escalators"*
- *"When riding with a pet, carry it"*

SYNTAX: GRAMMARS, PARSING & AMBIGUITY

- The parser search space is huge as for the effect of several forms of ambiguity that interacts in a combinatorial way

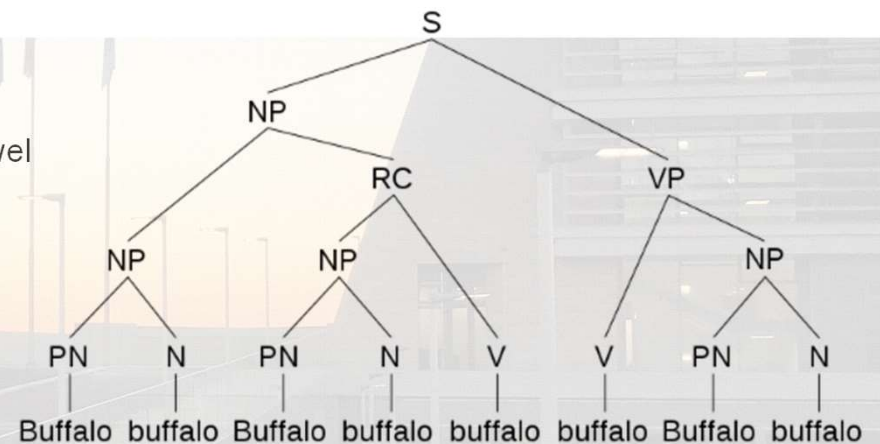
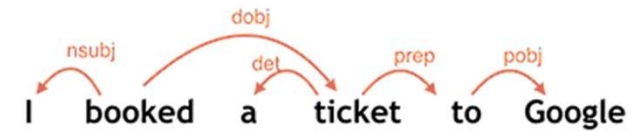
 e.g. *La vecchia porta la sbarra,*
 or *Buffalo buffalo Buffalo buffalo
buffalo buffalo Buffalo buffalo*

- Notice the strong relationship with semantics
 - Most of the ambiguities cannot be solved just at syntactic level
 - Lexical information (e.g. word senses) are crucial:

 *To operate in a market* viz. *To operate a body part*

 *Operare in un mercato* ≠ *Operare un paziente*

Dependency Parsing



Bison from Buffalo, New York who are intimidated by other bison in their community also happen to intimidate other bison in their community



(A(SHIP SHIPPING)SHIP) SHIPPING(SHIPPING SHIPS))

SEMANTICS



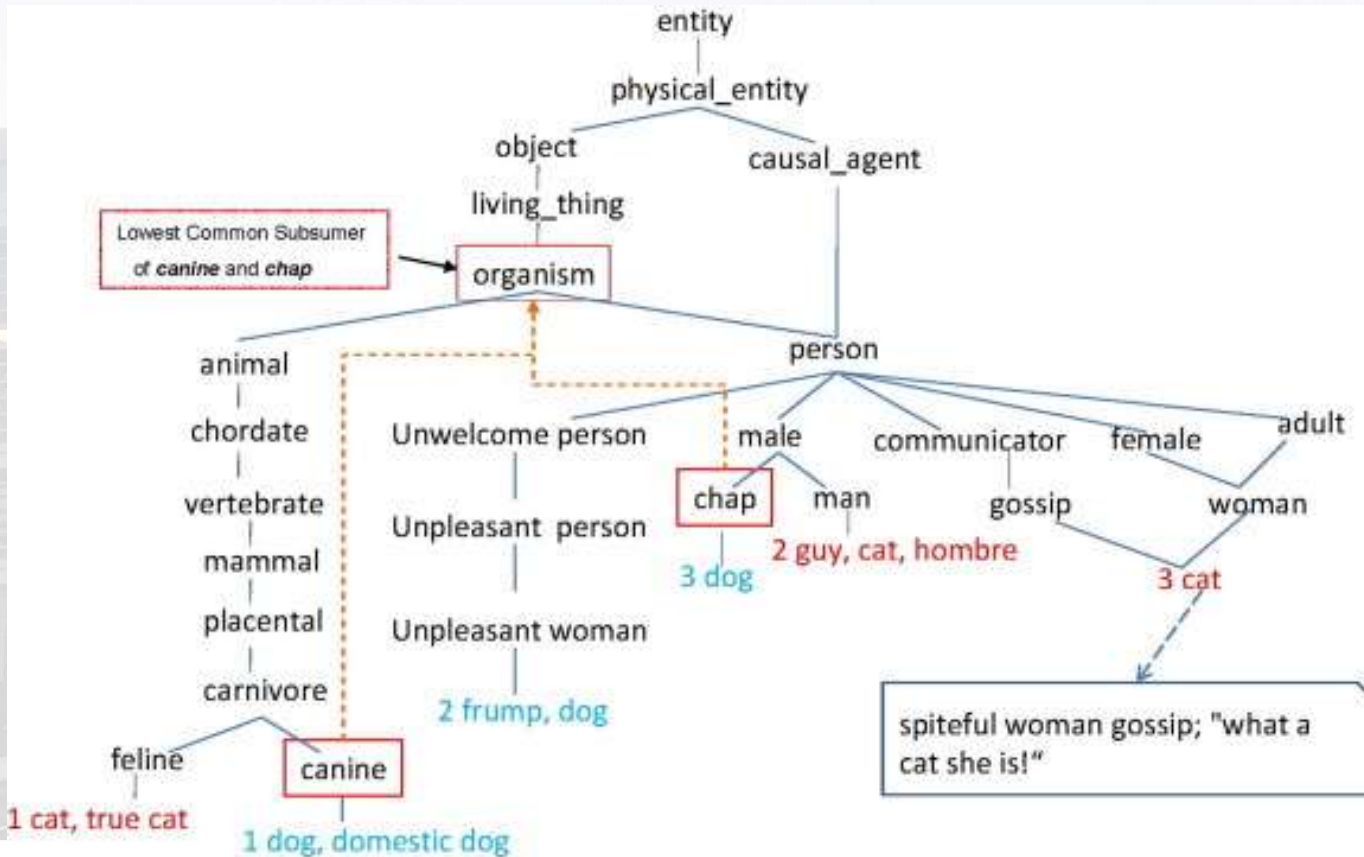
- What is the meaning of the sentence

John saw Kim?

- Desirable Properties:

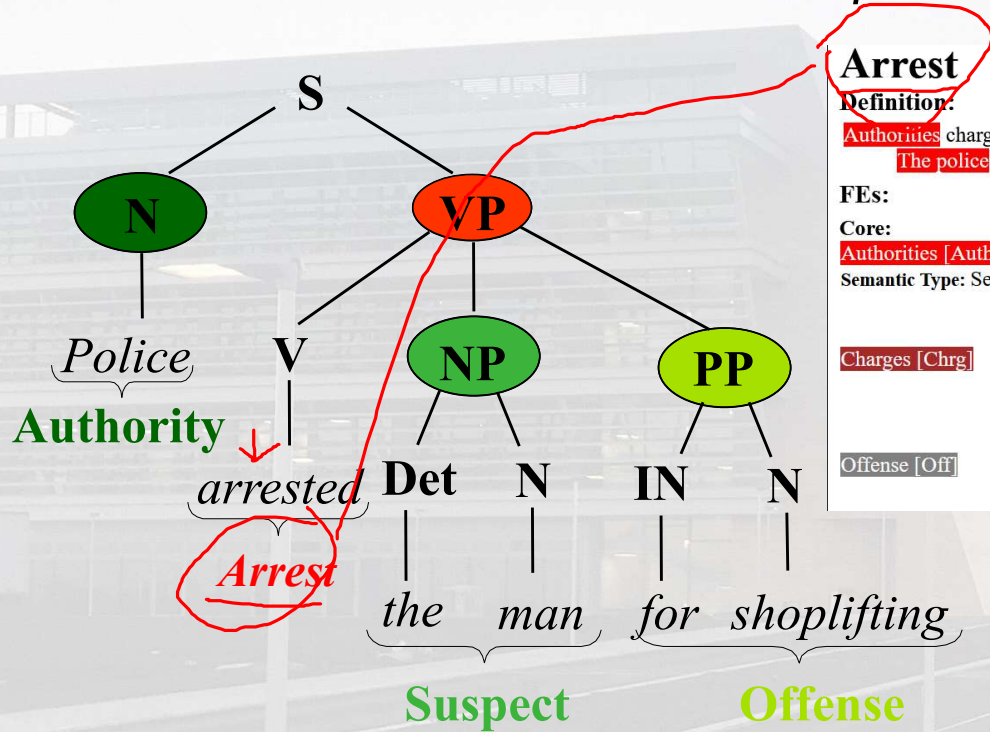
- It should be derivable as a function of the individual constituents, i.e. the meanings of constituents such as *Kim*, *John* and *see*
- Independent from syntactic phenomena, e.g. *Kim was seen by John* is a paraphrasis with the *same meaning*
- It must be directly used to trigger some inferences:
 - *Who was seen by John? Kim!*
 - *John saw Kim. He started running to her.*

WORD SENSES: THE WORDNET MODEL



FRAMENET: LINKING SYNTAX TO SEMANTICS

- Police arrested the man for shoplifting



Arrest

Definition:

Authorities charge a Suspect, who is under suspicion of having committed a crime (the Charges), and take him/her into custody.

The police ARRESTED Harry on charges of manslaughter.

FEs:

Core:

Authorities [Auth]

Semantic Type: Sentient

The Authorities charge the Suspect with committing a crime, and take him/her into custody.

The police ARRESTED Harry on charges of manslaughter.

Charges [Chrg]

Charges identifies a category within the legal system; it is the crime with which the Suspect is charged.

The police ARRESTED Harry on charges of manslaughter.

Offense [Off]

Offense identifies the ordinary language use of the reason for which a Suspect is arrested.

They arrested Harry for shoplifting.



FRAMENET LABELING: THE RELATIONAL VISION

Word	Predicate	Semantic Role ₁	Semantic Role ₂
Police	-	AUTHORITY	-
<i>arrested</i>	Target ₁	Arrest	-
the	-	SUSPECT	-
man	-	SUSPECT	-
for	-	OFFENSE	-
<i>shoplifting</i>	Target ₂	OFFENSE	Theft
merchandise	-	OFFENSE	GOODS

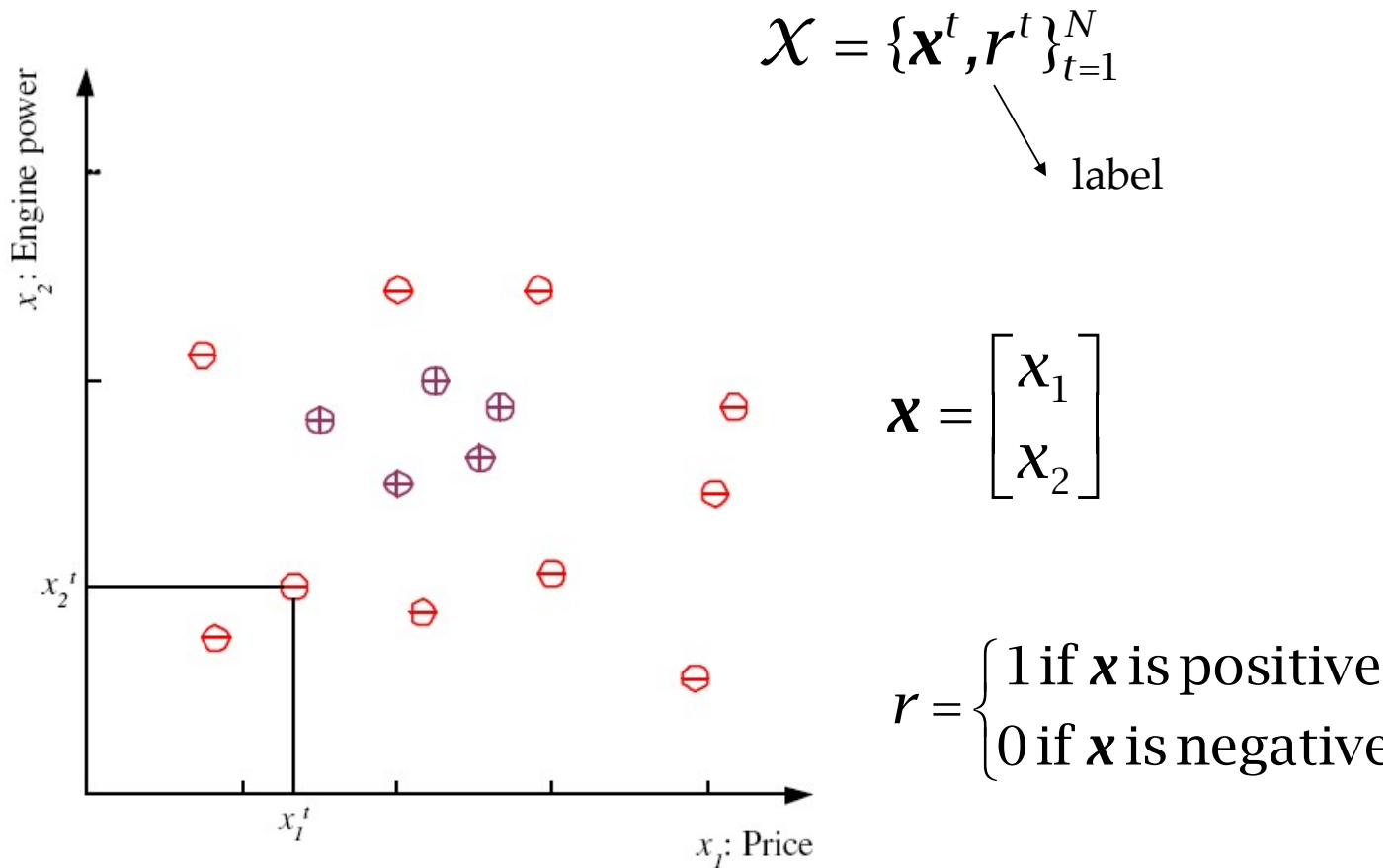


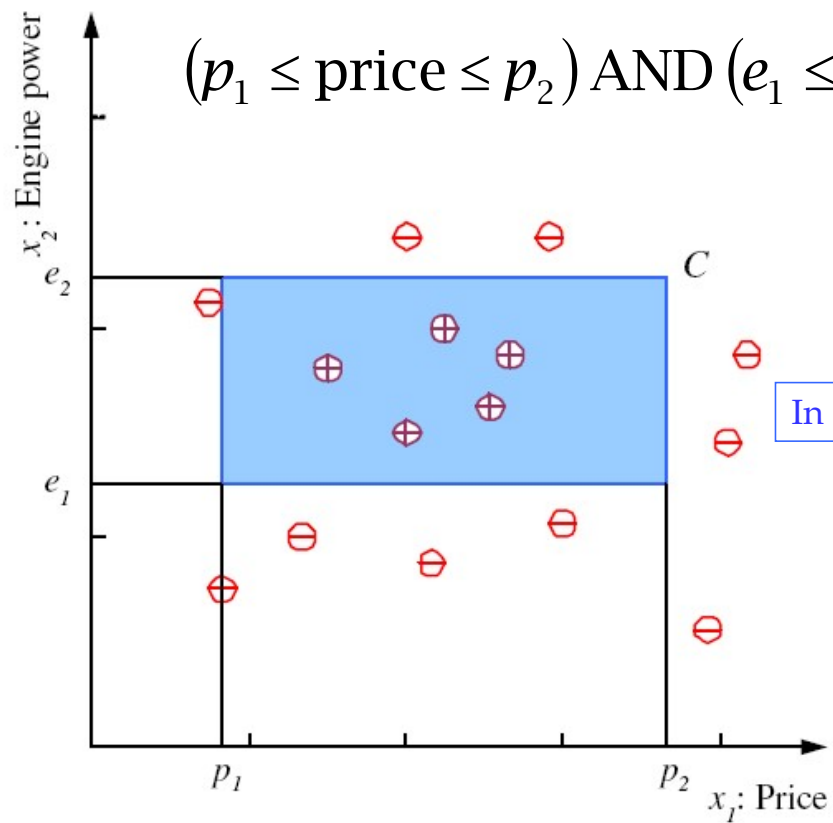
FROM STATISTICAL LEARNING THEORY TO SVMs

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

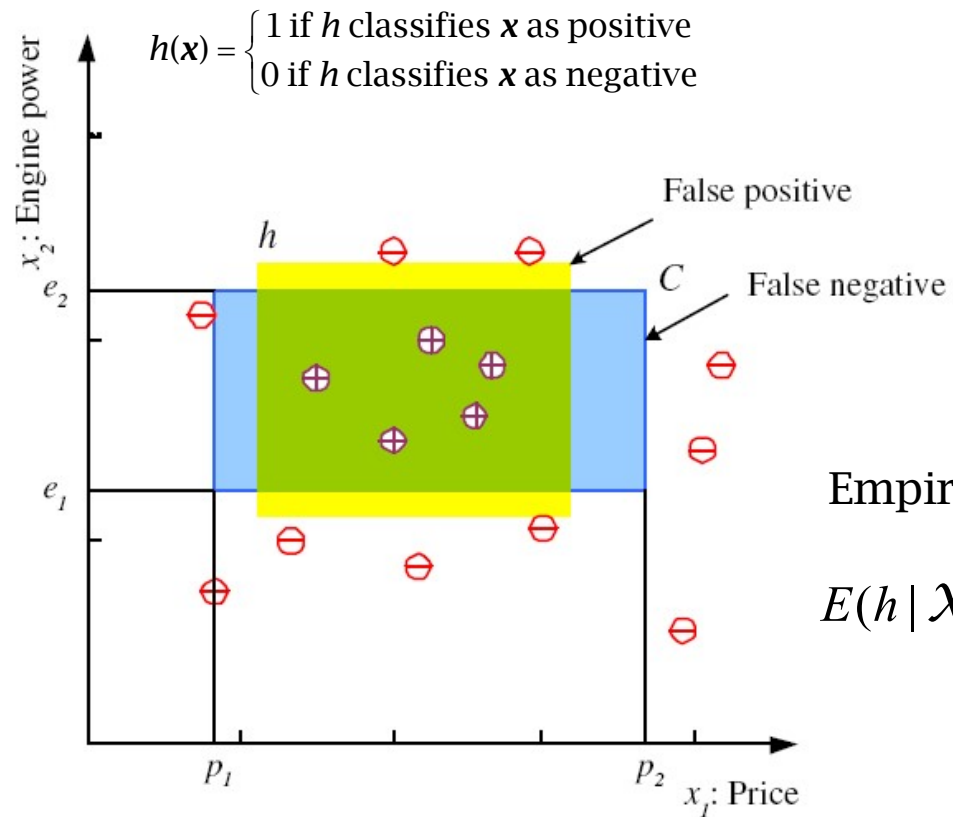
TRAINING SET \mathcal{X}



CLASS C 

In general we do not know $C(x)$.

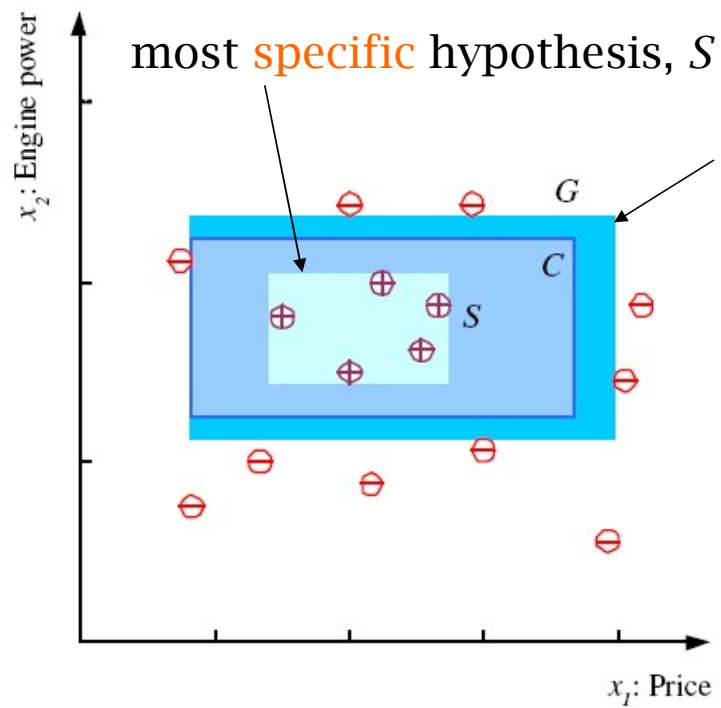
HYPOTHESIS CLASS \mathcal{H}



Empirical error:

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

S, G, AND THE VERSION SPACE





most **general** hypothesis, G

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

(Mitchell, 1997)

PROBABLY APPROXIMATELY CORRECT (PAC) LEARNING

- How many training examples are needed so that the tightest rectangle S which will constitute our hypothesis, will **probably** be **approximately correct**?
 - We want to be **confident** (*above a level*) that 
 - ... the **error probability is bounded** by some value 

- A concept class C is called **PAC-learnable** if there exists a PAC-learning algorithm such that, for any $\epsilon > 0$ and $\delta > 0$, there exists a fixed sample size such that, for any concept $c \in C$ and for any probability distribution on X , the learning algorithm produces a probably-approximately-correct hypothesis h
- a (PAC) **probably-approximately-correct hypothesis** h is one that has error at most ϵ with probability at least $1-\delta$.

PROBABLY APPROXIMATELY CORRECT (PAC) LEARNING

- In PAC learning, given a class C and examples drawn from some unknown but fixed distribution $p(x)$, we want to find the number of examples N , such that with probability at least $1-\delta$, h has error at most ε ? (Blumer et al., 1989)

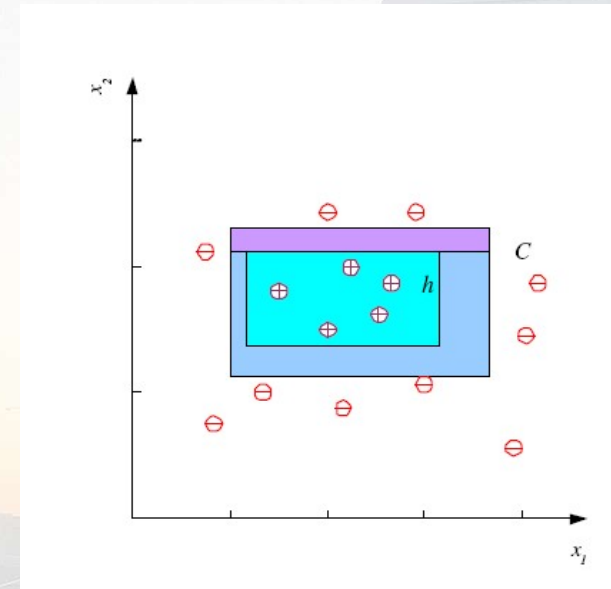
- $$P(C\Delta h \leq \varepsilon) \geq 1-\delta$$

- where $C\Delta h$ is (area of the) “the region of difference between C and h ”, and $\delta>0, \varepsilon>0$.

PAC LEARNING

How many training examples m should we have, such that with probability at least $1 - \delta$, h has error at most ϵ ? (Blumer et al., 1989)

- Let prob. of a + ex. in each strip be at most $\epsilon/4$
- Pr that a random ex. misses a strip: $1 - \epsilon/4$
- Pr that m random instances miss a strip:
 $(1 - \epsilon/4)^m$
- Pr that m random instances miss 4 strips:
 $4(1 - \epsilon/4)^m$
- We want $1 - 4(1 - \epsilon/4)^m \geq 1 - \delta$ or $4(1 - \epsilon/4)^m \leq \delta$
- Using $1 - x \leq e^{-x}$ an even stronger condition is:
 $[(1 - \epsilon/4)^m \leq \exp(-\epsilon m/4) \text{ so } (1 - \epsilon/4)^m \leq \exp(-\epsilon m/4) = \exp(-\epsilon m/4)]$
 $4e^{-\epsilon m/4} \leq \delta$ OR
- Divide by 4, take $\ln \dots$ and show that $m \geq (4/\epsilon) \ln(4/\delta)$



PROBABLY APPROXIMATELY CORRECT (PAC) LEARNING

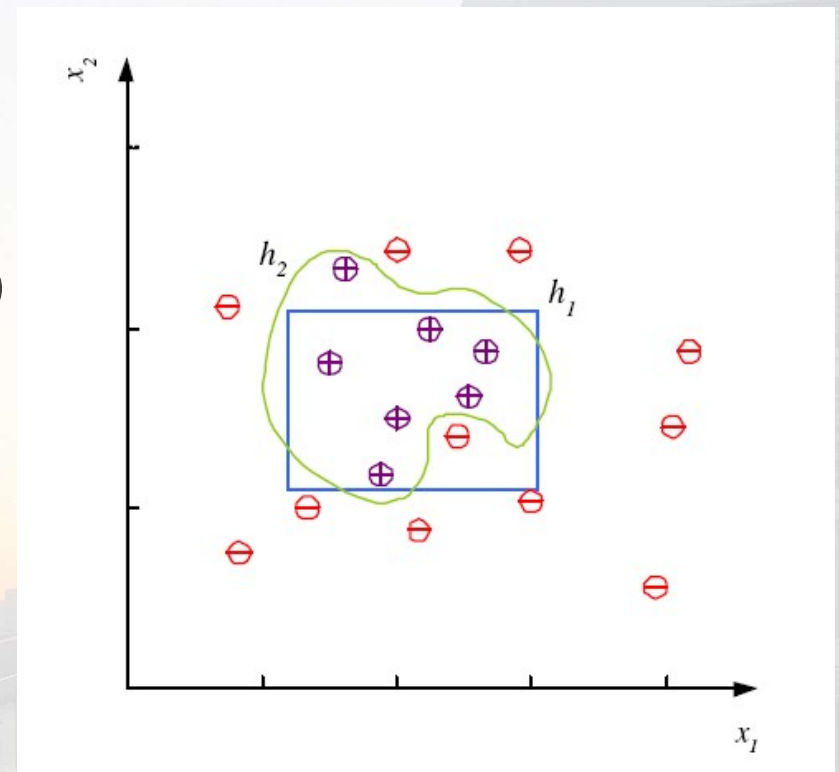
How many training examples m should we have, such that with probability at least $1 - \delta$, our hypothesis h has error at most ϵ ? (Blumer et al., 1989)

$$m \geq (4/\epsilon) \ln(4/\delta)$$

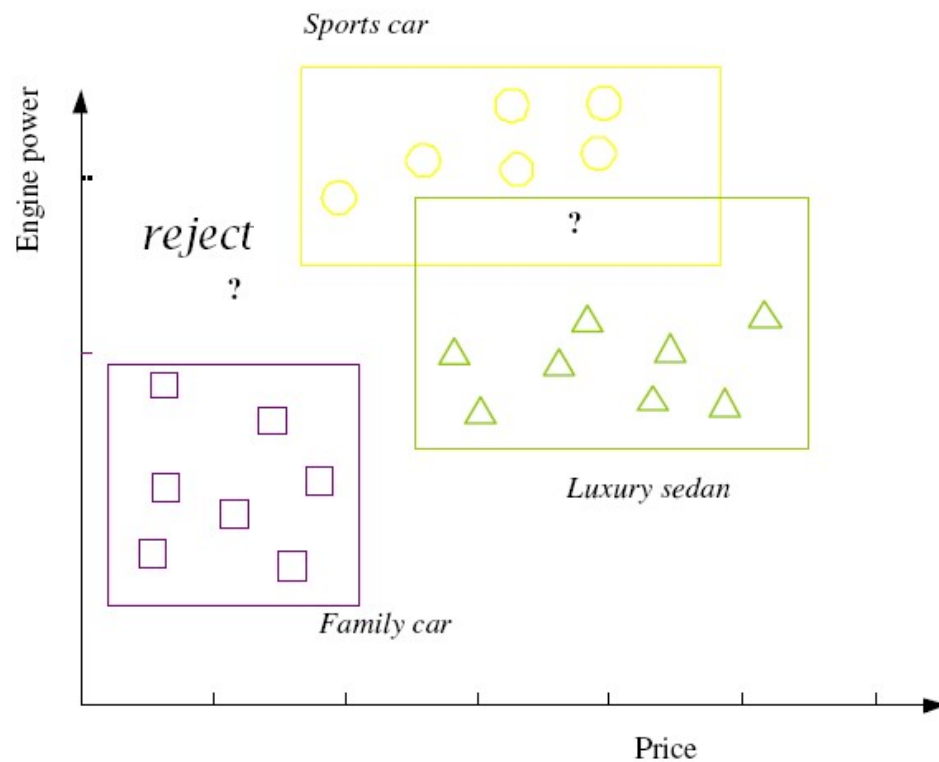
- m increases slowly with $1/\epsilon$ and $1/\delta$
- Say $\epsilon=1\%$ with confidence 95%, pick $m \geq 1752$
- Say $\epsilon=10\%$ with confidence 95%, pick $m \geq 175$

MODEL COMPLEXITY VS. NOISE

- 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, $C_i, i=1, \dots, K$



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

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

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

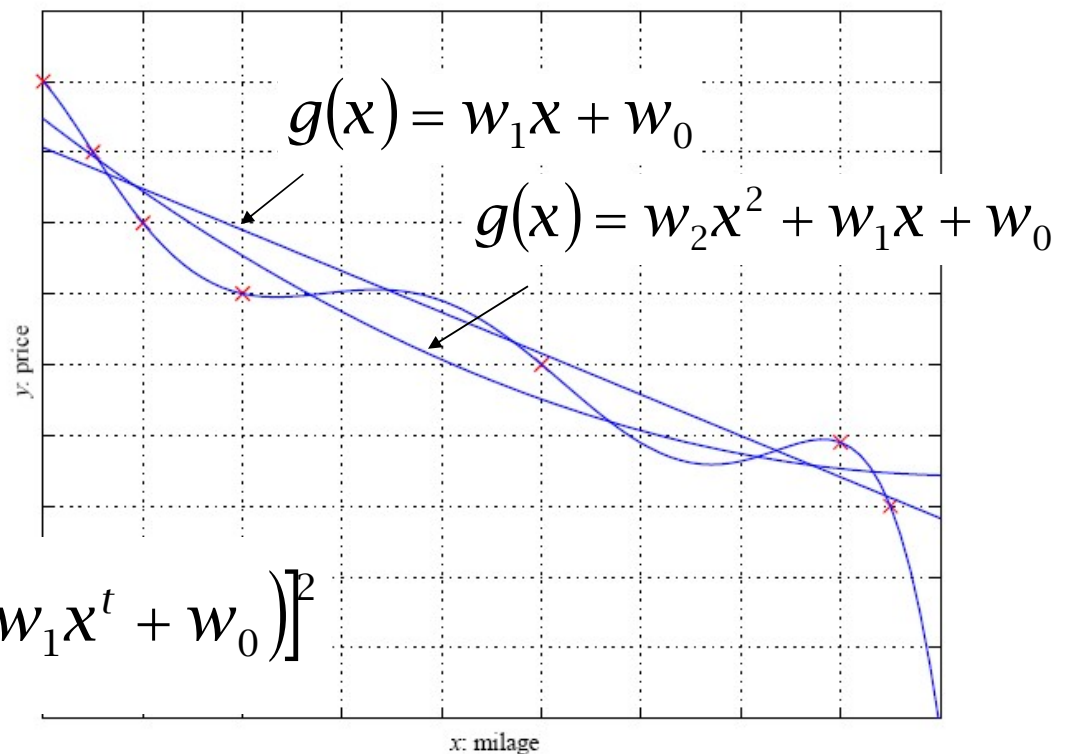
$$\mathcal{X} = \{x^t, r^t\}_{t=1}^N$$

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

$$r^t = f(x^t) + \varepsilon$$

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

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



VC (VAPNIK-CHERVONENKIS) DIMENSION

- N points can be labeled in 2^N ways as +/-
- \mathcal{H} **shatters** N if **there exists** a set of N points such that $h \in \mathcal{H}$ is consistent with **all** of these possible labels:
 - Denoted as: $VC(\mathcal{H}) = N$
 - Measures the capacity of H
- Any learning problem definable by N examples can be learned with no error by a hypothesis drawn from H

What is the VC dimension of axis-aligned rectangles?

FORMAL DEFINITION

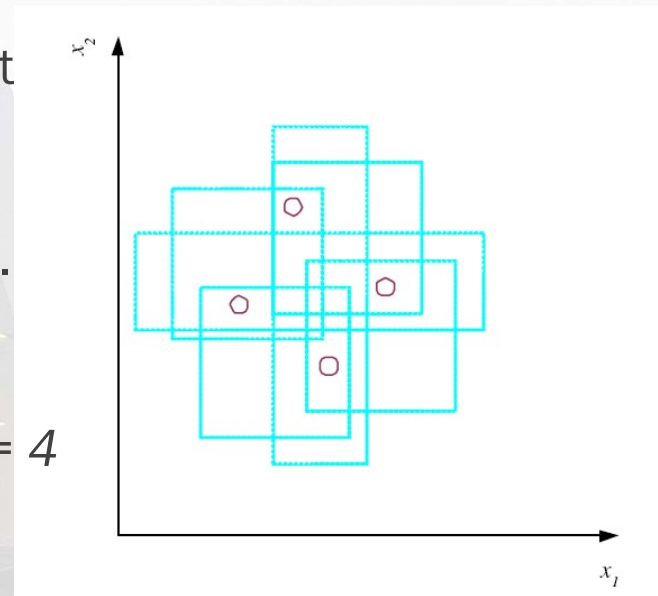
The VC Dimension

Definition: the VC dimension of a set of functions $H = \{h(\mathbf{x}, \alpha)\}$ is d if and only if there exists a set of points $\{x^i\}_{i=1}^d$ such that these points can be labeled in all 2^d possible configurations, and for each labeling, a member of set H can be found which correctly assigns those labels, but that no set $\{x^i\}_{i=1}^q$ exists where $q > d$ satisfying this property.

VC (VAPNIK-CHERVONENKIS) DIMENSION

- \mathcal{H} shatters N if there exists N points and $h \in \mathcal{H}$ such that h is consistent for any labelings of those N points.

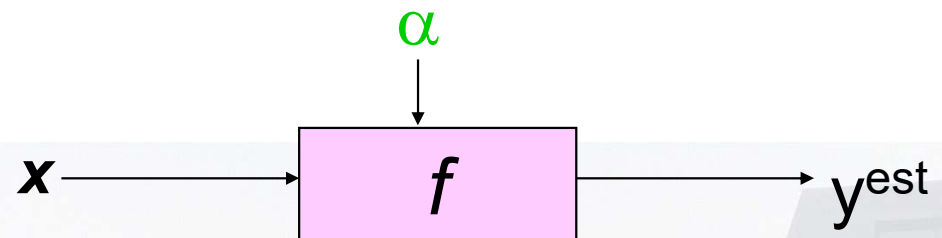
- $VC(\text{axis aligned rectangles}) = 4$



VC (VAPNIK-CHEVONENKIS) DIMENSION

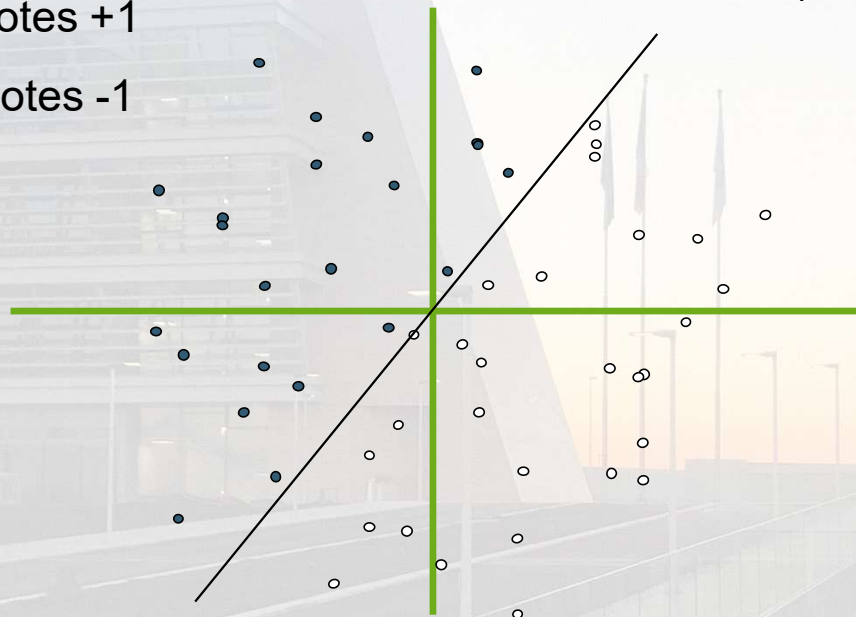
- *What does this say about using rectangles as our hypothesis class?*
- VC dimension is **pessimistic**: in general we do not need to worry about **all** possible labelings
- It is important to remember that one can choose the arrangement of points in the space, but then the hypothesis must be consistent with all possible labelings of those fixed points.

EXAMPLES

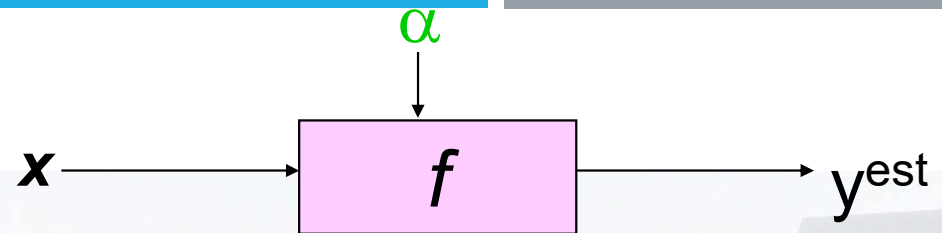


$$f(x, w) = \text{sign}(x \cdot w)$$

- denotes +1
- denotes -1

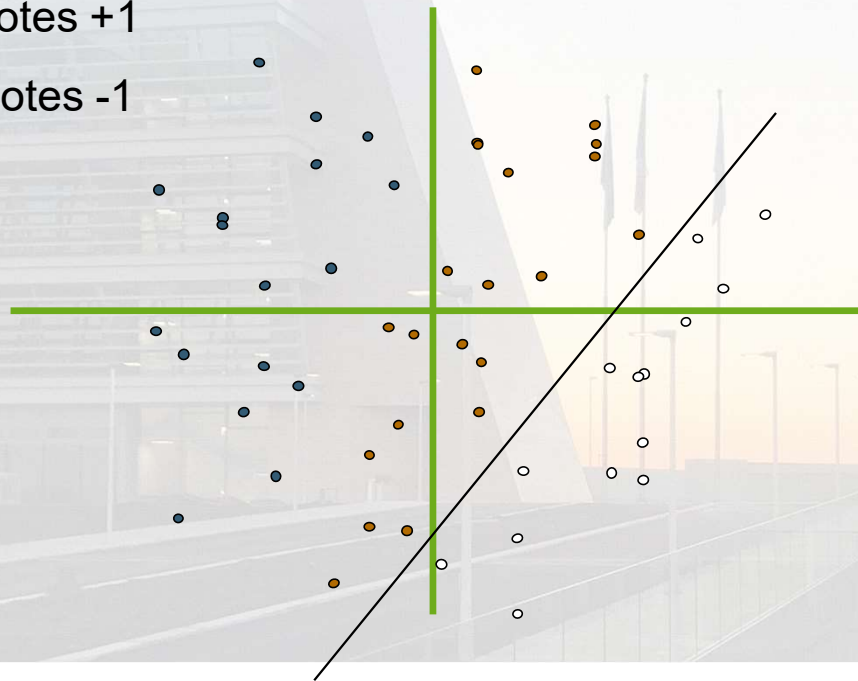


EXAMPLES



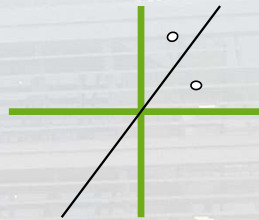
$$f(x, w, b) = \text{sign}(x \cdot w + b)$$

- denotes +1
- denotes -1



SHATTERING

- Question: Can the following f shatter the following points?

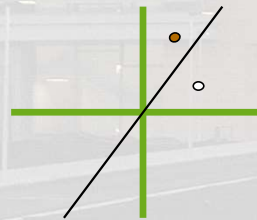


$$f(x, w) = \text{sign}(x \cdot w)$$

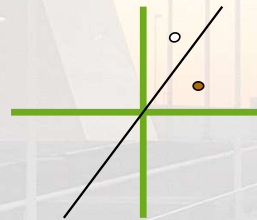
Answer: Yes. There are four possible training set types to consider:



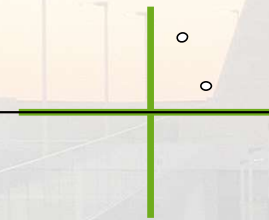
$$w = (0, 1)$$



$$w = (-2, 3)$$



$$w = (2, -3)$$



$$w = (0, -1)$$

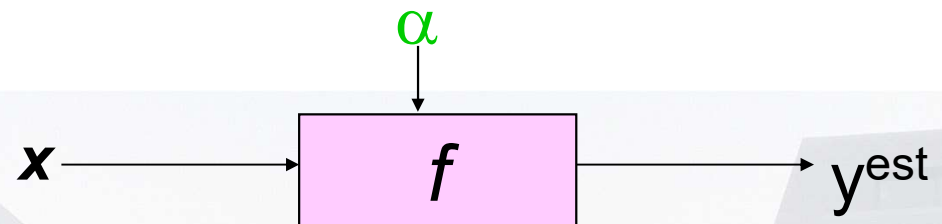
VC DIM OF LINEAR CLASSIFIERS IN M-DIMENSIONS

If input space is *m-dimensional* and if f is $\text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$, what is the VC-dimension?

$$h = m + 1$$

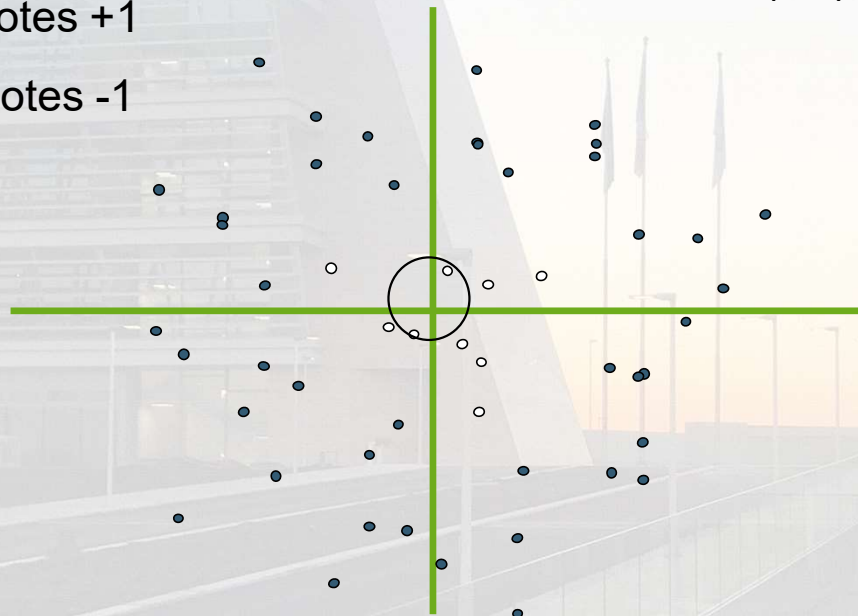
- Lines in 2D can shatter 3 points
- Planes in 3D space can shatter 4 points
- ...

EXAMPLES



$$f(x, b) = \text{sign}(x \cdot x - b)$$

- denotes +1
- denotes -1



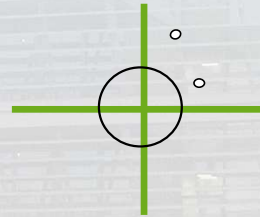
Diapositiva 38

rb1

roberto basili; 20/03/2023

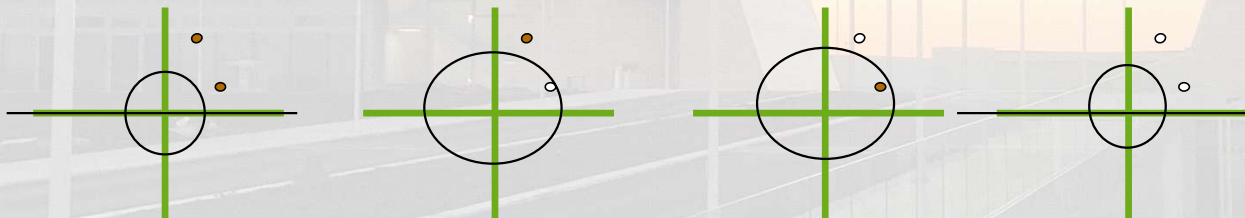
SHATTERING

- Question: Can the following f shatter the following points?



$$f(x, b) = \text{sign}(x \cdot x - b)$$

Answer: Yes. Hence, the VC dimension of circles on the origin is at least 2.



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 \mathcal{H}
- **Generalization**: How well a model performs on new data
- Different machines have different amounts of “power”.

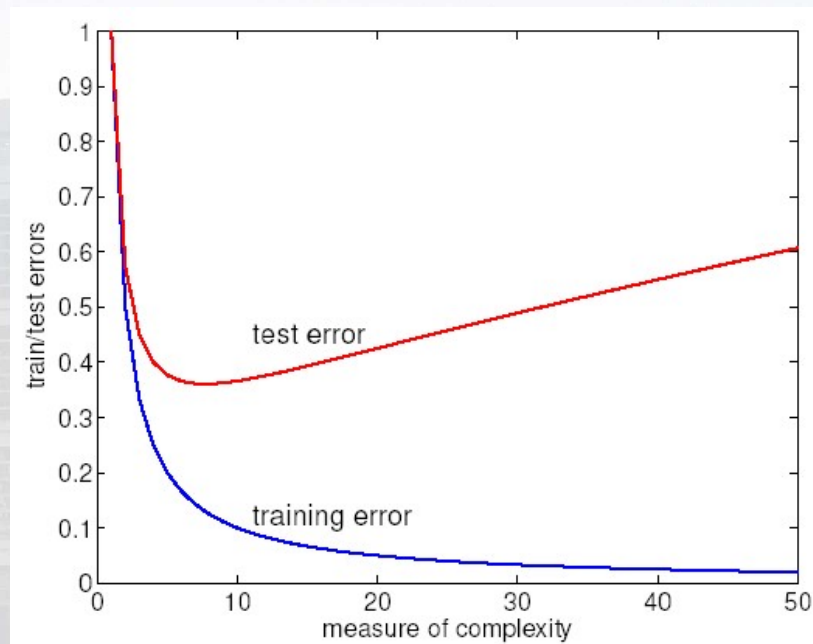
Tradeoff between:

- More power: Can model more complex classifiers but might overfit.
- Less power: Not going to overfit, but restricted in what it can model.
- **Overfitting**: \mathcal{H} more complex than C or f
- **Underfitting**: \mathcal{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 $N \uparrow$, $E \downarrow$
- As $c(\mathcal{H}) \uparrow$, first $E \downarrow$ and then $E \uparrow$

WHY CARE ABOUT COMPLEXITY?



- A quantitative measure of complexity is useful to determine the relationship between the training error (that we can observe during training) and the test error (which we want to minimize)

COMPLEXITY

- “Complexity” is a measure of a family of classifiers, not of any specific (fixed) classifier
- There are many possible measures for complexity
 - degrees of freedom (e.g. number of parameters in polynomials)
 - description length
 - Vapnik-Chervonenkis (VC) dimension
 - etc.

EXPECTED AND EMPIRICAL ERROR

$$\hat{\mathcal{E}}_n(i) = \frac{1}{n} \sum_{t=1}^n \overbrace{\text{Loss}(y_t, h_i(\mathbf{x}_t))}^{=0,1} = \text{empirical error of } h_i(\mathbf{x})$$
$$\mathcal{E}(i) = E_{(\mathbf{x}, y) \sim P} \{ \text{Loss}(y, h_i(\mathbf{x})) \} = \text{expected error of } h_i(\mathbf{x})$$

LEARNING AND THE VC DIMENSION

- Let d_{VC} be the VC-dimension of our set of classifiers F .

Theorem: With probability at least $1 - \delta$ over the choice of the training set, for all $h \in F$

$$\mathcal{E}(h) \leq \hat{\mathcal{E}}_n(h) + \epsilon(n, d_{VC}, \delta)$$

where

$$\epsilon(n, d_{VC}, \delta) = \sqrt{\frac{d_{VC}(\log(2n/d_{VC}) + 1) + \log(1/(4\delta))}{n}}$$

MODEL SELECTION

- We try to find the model with the best balance of complexity and the fit to the training data
- Ideally, we would select a model from a nested sequence of models of increasing complexity (VC-dimension)

Model 1 d_1

Model 2 d_2

Model 3 d_3

where $d_1 \leq d_2 \leq d_3 \leq \dots$

- The model selection criterion is: find the model class that achieves the lowest upper *bound* on the expected loss

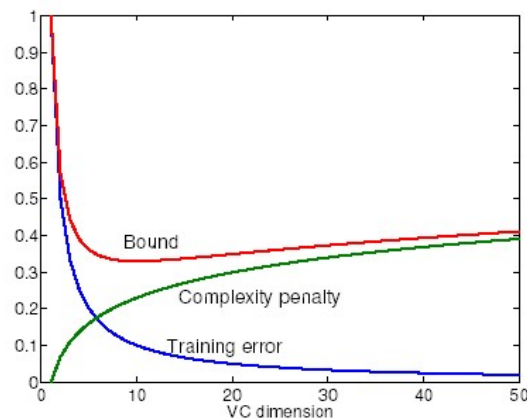
Expected error \leq Training error + Complexity penalty

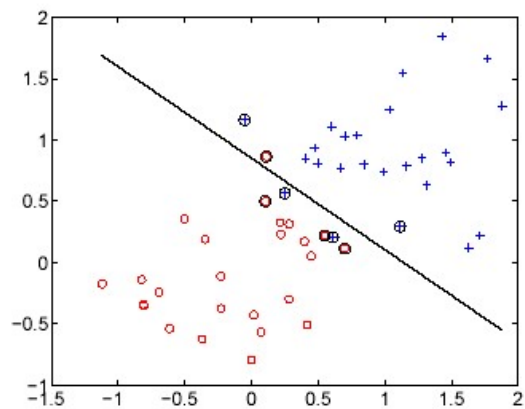
VC DIMENSION AND STRUCTURAL RISK MINIMIZATION

- We choose the model class F_i that minimizes the upper bound on the expected error:

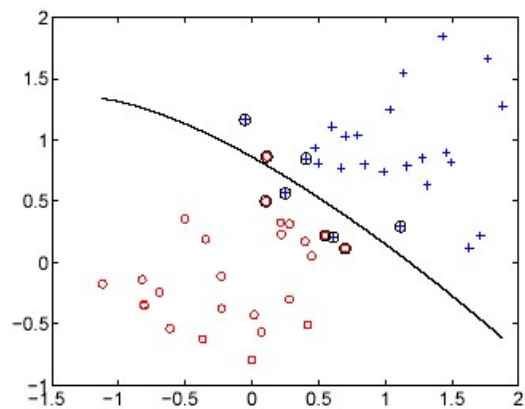
$$\mathcal{E}(\hat{h}_i) \leq \hat{\mathcal{E}}_n(\hat{h}_i) + \sqrt{\frac{d_i(\log(2n/d_i) + 1) + \log(1/(4\delta))}{n}}$$

where \hat{h}_i is the best classifier from F_i selected on the basis of the training set.

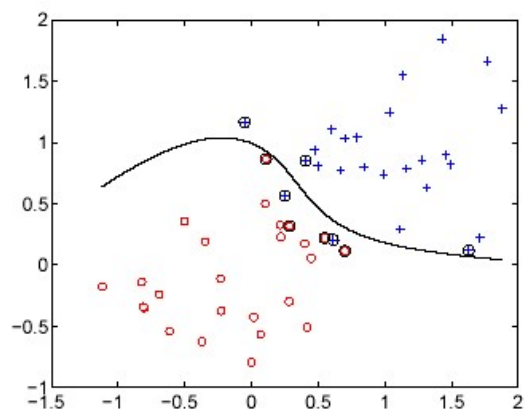




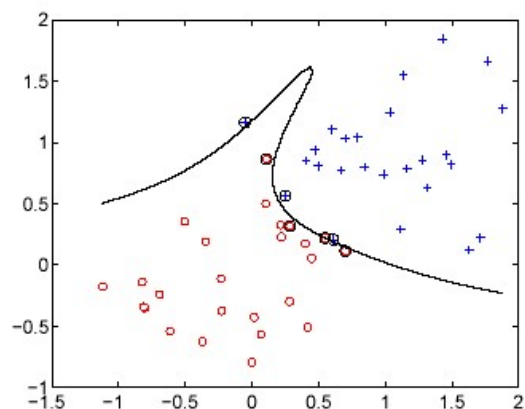
linear



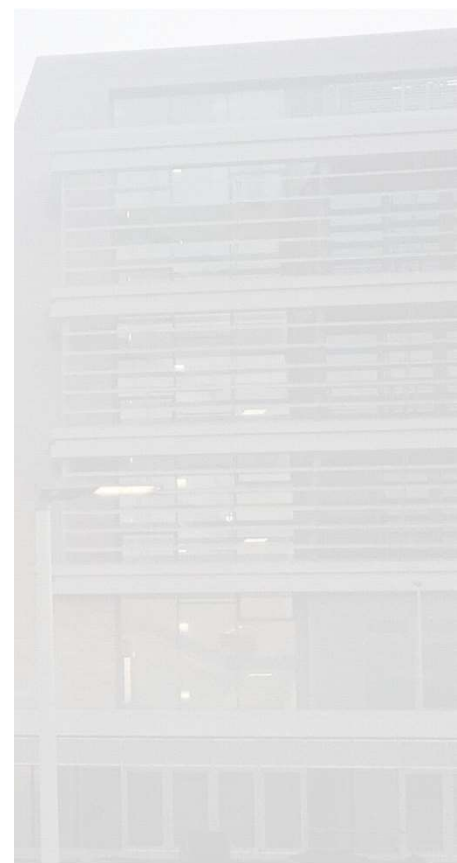
2nd order polynomial



4th order polynomial



8th order polynomial



STRUCTURAL RISK MINIMIZATION

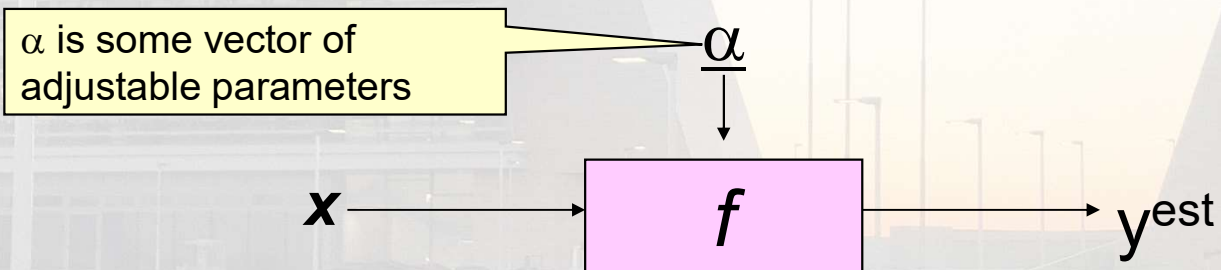
- Number of training examples $n = 50$, confidence parameter $\delta = 0.05$.

Model	d_{VC}	Empirical fit	$\epsilon(n, d_{VC}, \delta)$
1 st order	3	0.06	0.5501
2 nd order	6	0.06	0.6999
4 th order	15	0.04	0.9494
8 th order	45	0.02	1.2849

- Structural risk minimization would select the simplest (linear) model in this case.

SUMMARY: A LEARNING MACHINE

- A learning machine f takes an input x and transforms it, somehow using factors (as weights) $\underline{\alpha}$, into a predicted output $y^{est} = +/- 1$



VC-DIMENSION AS MEASURE OF COMPLEXITY

$$\text{TESTERR}(\vec{\alpha}) \leq \text{TRAINERR}(\vec{\alpha}) + \sqrt{\frac{h(\log(2R/h) + 1) - \log(\eta/4)}{R}}$$

i	f_i	TRAINERR	VC-Conf	Probable upper bound on TESTERR	Choice
1	f_1				
2	f_2				
3	f_3				?
4	f_4				
5	f_5				
6	f_6				

USING VC-DIMENSIONALITY

- People have worked hard to find VC-dimension for ...
 - Decision Trees
 - Perceptrons
 - Neural Nets
 - Decision Lists
 - Support Vector Machines
 - ...and many many more
- All with the goals of
 - Understanding which learning machines are more or less powerful under which circumstances
 - Using Structural Risk Minimization for to choose the best learning machine






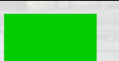


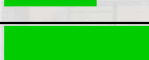




ALTERNATIVES TO VC-DIM-BASED MODEL SELECTION

Cross Validation

- To estimate generalization error, we need data unseen during training. We split the data as:
 - Training set (50%) $M1$ $M2$ $\text{train}(M2) < \text{train}(M1)$
 - Validation set (25%) $\text{test}(M1, V_s) = P1$ $\text{test}(M2, V_s) = P2$ $P2 > P1$
 - Test (publication) set (25%)
- Resampling when there is few data
 - N-fold cross-validation: N-2 fold for training, 1 fold as validation set and 1 fold for testing ($N \cdot (N-1)$ tests)

ALTERNATIVES TO VC-DIM-BASED MODEL SELECTION

- What could we do instead of the scheme below?
Cross-validation

i	f_i	TRAINER R	10-FOLD-CV-ERR	Choice
1	f_1			
2	f_2			
3	f_3			
4	f_4			
5	f_5			
6	f_6			

EXTRA COMMENTS

- An excellent tutorial on VC-dimension and Support Vector Machines

C.J.C. Burges. A tutorial on support vector machines for pattern recognition. *Data Mining and Knowledge Discovery*, 2(2):955-974, 1998.

WHAT YOU SHOULD KNOW

- Definition of PAC learning
- The definition of a learning machine: $f(x, \alpha)$
- The definition of Shattering
- Be able to work through simple examples of shattering
- The definition of VC-dimension
- Be able to work through simple examples of VC-dimension
- Structural Risk Minimization for model selection
- Awareness of other model selection methods