Malware Detection with Machine Learning Methods (IN2207)

Thomas Stibor

Fakultät für Informatik Technische Universität München

thomas.stibor@in.tum.de

SS2011

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Organizational

- Language: English
- Lecture: Tuesday, 13:00 14:30, MI 01.06.011
- Exercise course: Friday, 10:00 10:45, MI 01.06.011
- Credits: 2+1 SWS / 4,0 ECTS
- Final exam (oral)

The homeworks will contain written questions and questions that require some R programming.

Exercise course on Friday 06.05.2011 will be an introduction to R (programming).

Overview

- Theoretical Limits of Virus Scanners
- Brief Warm-Up in Probability
- Naive Bayes Classifier
- Probabilistic Models (HMM, Topic Models)
- Linear Classifiers
- Neural-Networks and Support Vector Machines
- Feature Selection & Extraction
- Anomaly Detection
- Adversarial Learning

Literature



Some figures are taken from Bishop's new book (Pattern Recognition and Machine Learning).

Theoretical Limits of Virus Scanners

Can any virus be detected? What are the limits?

 Identifying a mutating virus is NP-complete.
 Reliable Identification of Bounded-length Viruses is NP-complete, Diomidis Spinellis, IEEE Transactions on Information Theory, Vol. 49, No. 1 (2003), pp. 280-284.

• Detecting viruses is undecidable.

Computer Viruses: Theory and Experiments, Fred Cohen, Computers & Security, 6 (1987), pp. 22-35.

Computer Virus

Definition: A computer virus is a program that can infect other programs by modifying them to include a possibly evolved copy of it self. [Fred Cohen, Computer Viruses: Theory and Experiments, 1987]

```
program virus := 1234567;
  subroutine infect-executable := {
    loop:file = get-random-executable-file;
    if first-line-of-file = 1234567 then goto loop:
    prepend virus to file;
  3
  subroutine do-damage := {whatever damage is to be done}
  subroutine trigger-pulled := {return true if some condition}
  main-program := { infect-executable;
    if trigger-pulled then do-damage;
    goto next;
  }
  next: /* execute rest of the program it was prepended to. */
}
```

(B)

Identifying a Virus is $\mathcal{NP}\text{-}\mathsf{complete}$

Idea:

Show that a virus detector D for a certain virus strain V can be used to solve the satisfiability problem, which is known to be \mathcal{NP} -complete.

- The satisfiability of the problem which is examined is however *not* a special case.
- The virus V is a mutating self-replicating program.
- Assume that the virus detector *D* can reliably determine in polynomial time whether a given candidate program *C* is a mutation of the virus *V*.
- We will use the virus detector as an oracle for determining the satisfiability of an *N*-term Boolean formula *S*.

Link between SAT and Virus Instance

Use the satisfiability formula S to create a virus archetype A and a possible instance of a virus phenotype P. The virus is a triple

where

- f is the virus processing and replication function,
- *s* is a Boolean value indicating whether an instance of the virus has found a solution to *S*,
- *c* is an integer encoding the candidate values for *S*.

The function f maps a triple (f, s, c) into a new triple (f, s', c') and is defined as follows

$$\lambda(f, s, c).(f, s \lor S, \text{ if } c = 2^N \text{ then } c \text{ else } c + 1).$$

Link between SAT and Virus Instance (cont.)

Each term x_n in S is calculated from c through the expression

$$\left\lfloor \frac{c}{2^n} \right\rfloor \mod 2 \stackrel{?}{=} 1$$
 which results in a Boolean value.

bit position <i>c</i>			decimal value
2	1	0	п
0	0	0	0
0	0	1	1
0	1	0	2
0	1	1	3
1	0	0	4
1	0	1	5
1	1	0	6
1	1	1	7

Is bit at position c = 2 set to 1 for decimal value n = 4? Result is TRUE.

Is bit at position c = 0 set to 1 for decimal value n = 6? Result is FALSE.

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Will a Virus Mutation ever Satisfy S

A new generation of the virus is generated by applying f to the current generation.

- evaluates S by extracting successive Boolean value combinations from c,
- 2 increments c until it reaches 2^N ,
- \bigcirc passes the result of the S evaluation to the next generation.

We can now ask virus detector D whether the virus archetype A

(f, FALSE, 0)

will ever result in a virus mutation phenotype P

 $(f, \text{TRUE}, 2^N)$

that is whether one of the virus mutations will satisfy S.

Example

Consider the satisfiability of the formula

$$S = (x_0 \vee x_1) \wedge \overline{x}_0.$$

The virus replication function f will be

$$\lambda(f, s, c).(f, s \lor (x_0 \lor x_1) \land \overline{x}_0, c+1)$$

the corresponding archetype A

$$(\lambda(f,s,c).(f,s\vee(x_0\vee x_1)\wedge\overline{x}_0,c+1),F,0),$$

where F denotes FALSE and T denotes TRUE.

A B K A B K

Example (cont.)

This particular virus will generate a mutation P and thereby indicate that S is satisfiable in four generations through the following sequence

$$\begin{split} & (\lambda(f,s,c).(f,s\vee S,c+1), \mathsf{F}\vee(\mathsf{F}\vee\mathsf{F})\wedge\overline{\mathsf{F}},0+1) \quad \rightarrow \\ & (\lambda(f,s,c).(f,s\vee S,c+1),\mathsf{F},1) \quad \rightarrow \\ & (\lambda(f,s,c).(f,s\vee S,c+1),\mathsf{F}\vee(\mathsf{T}\vee\mathsf{F})\wedge\overline{\mathsf{T}},1+1) \quad \rightarrow \\ & (\lambda(f,s,c).(f,s\vee S,c+1),\mathsf{F},2) \quad \rightarrow \\ & (\lambda(f,s,c).(f,s\vee S,c+1),\mathsf{F}\vee(\mathsf{F}\vee\mathsf{T})\wedge\overline{\mathsf{F}},2+1) \quad \rightarrow \\ & (\lambda(f,s,c).(f,s\vee S,c+1),\mathsf{T}\vee(\mathsf{T}\vee\mathsf{T})\wedge\overline{\mathsf{T}},3+1) \quad \rightarrow \\ & (\lambda(f,s,c).(f,s\vee S,c+1),\mathsf{T}\vee(\mathsf{T}\vee\mathsf{T})\wedge\overline{\mathsf{T}},3+1) \quad \rightarrow \\ & (\lambda(f,s,c).(f,s\vee S,c+1),\mathsf{T},4) \quad \equiv \quad P. \end{split}$$

(4 間) トイヨト イヨト

Summary

It was shown that a reliable virus detector D operating in polynomial time can be used as a polynomial time satisfiable oracle and hence, reliable detection of a mutating virus is \mathcal{NP} -complete

TRANSACTEORS ON DEFORMATION THEORY, VOL. 19, NO. 1, IAMUARY 2003 (DRAFT

Reliable Identification of Bounded-length Viruses is NP-complete

Diomidis Spinellis, Member, IEEE

Advocve—A view is a program that replicates itself by opping its ords into other his. A common view protection mechanism inviews examing like is detect one patterns of known diverse. The prover that the produces of redshift shifting a bounded length prover that the produces of redshift shifting a bounded length a cartian views itemis can be used it a solver the autholibility prodmons. The implication of this result is observe the autholibility prodable of the body generation gates as views mattalianding and advocation of the result is observed in an other solve the device internet and a which devices and arrows result.

Infer Terms-buffer-overflow, complexity, detection, identification, mutation, NP-complete, security, virus

I. INTRODUCTION

Other three sound defines against compare remains the new sound series of an annual series program that do also that also been sound to also the sound of the sound series and th

II. VIRAL SOFTWARE

Interiorally created maticass offware [2]--uclea termed and-non-over, typically classified that Trajka herese, visues, and versus [3]. A Trajka heres is a program that caplain the rights of its sour to perform an action is user down oro inend, a visus is a Trajka hence that registants itself by copying its codes into other program that registants itself by copying its code sourch weakpowers to just down on compare.

D. Spitellis is an Assistant Prolosser in the Department of Management Science and Technology at the Athons University of Economics and Business, Athens, Gonce, Lemail: definition of prorident and the Athensister and Athensister and Athensister and Athensister and Proceedings of a working paper definition tail for a "This is a maximum conduction Discover and an analyzed and an and a second proparative conductive conductive conductive conduction paper definite and the loss to "This is a maximum conductive conductiv

These indexactions on the processing of the class cards during order. These is the classical of the classic

A matter of visus prevanitor and denotion methods how how preposed and are commonly implemented [15, [6]]. Reence cT/ constants at anomales Historyappi of malware mailtime and the second second second second second second and the second second second second second second second and the second second second second second second second many high-second second s

Therefore, as a secondary lise of definite, detection resears are offer enrophotod to locat virus instances and indexsonance in the enrophotod on locat virus instances and indexcomparison of the system's programs against latence good terstory larged bootstances in the system of the system's program in the system's programs against latence good tersing and system signature. By one disc comparison of the system's sector signature is the system sector sector of the system's regression and the system's sector sector and the system's sector signature is the system's sector sector sector sector sector sector sector sector sector of the system's sector sector

View series however, how developed a series of examplements, being only address camples of view down of develop sequenced is hirder that detacts at the stern p1 of the development of the development of the stern p1 performance of the development of the stern p1 performance at the development of the development of performance at the development of the development of a variable corporation of the development of the and advectory the method performance and the development of the development of the and the development of the development of the advectory and the development of the development of the advectory and the development of the development of the development of advectory the method of the outperformance and the development of the development of the development of advectory the development of the outperformance advectory the development of the outperformance advectory the development of advectory the development of the outperformance advectory the development of the outperformance

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Course IN2207

Virus Detection is Undecidable

Virus detection is undecidable, that is, there is no program that detects all viruses.

To proof this statement one has to show the undecidability of the language

 $V_{\text{TM}} = \{ \langle M, \langle P \rangle \rangle \mid M \text{ is a TM and } M \text{ detects whether } P \text{ is a virus} \}.$

Proof sketch:

Suppose that V_{TM} is decidable and H is a decider for V_{TM} .

$$H(\langle M, \langle P \rangle \rangle) = \begin{cases} \text{TRUE} & \text{if } M \text{ detects } P \text{ as a virus} \\ \text{FALSE} & \text{if } M \text{ fails to detect } P \text{ as a virus.} \end{cases}$$

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Virus Detection is Undecidable (cont.)

- On input (M, (P)), where M is a TM and (P) a string which encodes program P, H outputs TRUE if M detects P as a virus. In contrast, H outputs FALSE if M fails to detect P as a virus.
- Construct a new Turing machine D with H as a subroutine.
- This new TM calls H to determine what M does when the input to M is its own description (M).
- Once *D* has determined this information, it does the opposite. That is, it outputs FALSE if *M* outputs TRUE and it outputs TRUE if *M* outputs FALSE. The following is a description of *D*.

A B K A B K

Virus Detection is Undecidable (cont.)

- D =On input $\langle M \rangle$, where M is a TM:
 - **1** Run *H* on input $\langle M, \langle M \rangle \rangle$.
 - Output the opposite of what H outputs; that is, if H outputs TRUE, then output FALSE and if H outputs FALSE, then output TRUE.

More formally

$$D(\langle M \rangle) = \begin{cases} \text{TRUE} & \text{if } M \text{ fails to detect } \langle M \rangle \text{ as a virus} \\ \text{FALSE} & \text{if } M \text{ detects } \langle M \rangle \text{ as a virus.} \end{cases}$$

What happens when we run D with its own description $\langle D \rangle$ as input?

A B K A B K

Virus Detection is Undecidable (cont.)

In that case we get

$$D(\langle D \rangle) = \begin{cases} \text{TRUE} & \text{if } D \text{ fails to detect } \langle D \rangle \text{ as a virus} \\ \text{FALSE} & \text{if } D \text{ detects } \langle D \rangle \text{ as a virus.} \end{cases}$$

No matter what D does, it is forced to do the opposite and hence this results in a contraction.

In summary

- H outputs TRUE on input (M, (P)) exactly when M detects P as a virus.
- D fails to detect $\langle M\rangle$ as a virus exactly when M detects $\langle M\rangle$ as a virus.
- D fails to detect ⟨D⟩ as a virus exactly when D detects ⟨D⟩ as a virus.

Machine Learning for Information Security Problems

- Spam detection: Text analysis, document clustering.
- Intrusion detection: Anomaly detection, one-class learning.
- Detection of malicious behavior of processes: Sequential data analysis of system calls.
- Credit card fraud detection: Time series analysis of transactions.

Machine learning methods can be applied for such kind of problems.

The goal of this course is to learn about the machine learning methods.

o ...

A B A A B A

An Introductory Example

Suppose that a fishpacking factory wants to automate the process of sorting incoming fish (salmon and sea bass).



After some preprocessing, each fish is characterized by feature vector $\mathbf{x} = (x_1, x_2) \in \mathbb{R}^2$ (pattern), where the first component is the lightness and the second component the length.

T.Stibor (TUM)

Course IN2207

SS2011 19 / 256

Pattern belongs to Class?



Given labeled training data $(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_N, y_N) \in \mathbb{R}^n \times Y$ coming from some unknown probability distribution $P(\mathbf{x}, y)$. In this example, $Y \in \{\text{salmon, sea bass}\}$ and n = 2. Unseen (unlabeled) pattern belongs to class salmon or sea bass?

A (too underfitted) Classifier



This linear separation suggests the rule: Classify the fish as salmon if its features falls below the *decision boundary*, otherwise as sea bass.

A (too overfitted) Classifier



A too complex model will lead to decision boundary that gives perfect classification accuracy on trainig set (seen patterns), but poor classification on *unseen* patterns.

A Good Classifier



Optimal tradeoff between performance on the training set and simplicity of the model. This gives high classification accuracy on unseen patterns, i.e. it gives good *generalization*.

Probability Theory

- Uncertainty is key concept in machine learning.
- Probability provides consistent framework for the *quantification* and *manipulation* of uncertainty.

Probability of an event is the fraction of times that event occurs out of the total number of trials, in the limit that the total number of trials goes to infinity.

Elementary rules of probability

- Sum rule.
- Product rule.

Derive Sum and Product Rules



- Random variable X can take values {x_i}, i = 1, 2, ..., M (M = 5), Y can take values {y_j}, j = 1, 2, ..., L (L = 3).
- Consider a total number N of instances of these variables, and denote the number of instances where $X = x_i$ and $Y = y_i$ by n_{ij} (number of points in the corresponding cell of the array).

Derive Sum and Product Rules (cont.)



- Number of points in column *i*, corresponding to $X = x_i$, is denoted by c_i .
- Number of points in row *j*, corresponding to $Y = y_i$, is denoted by r_i .

Joint Probability

Probability that X will take value x_i and Y will take the value y_i is written

$$p(X = x_i, Y = y_i) = \frac{n_{ij}}{N}$$

where we are implicitly considering the limit $N \to \infty$. Put it the other way around:

• Number of points falling in the cell *i*, *j* as a fraction of the total number of points.

Probability that X takes the value x_i irrespective of the value of Y is

$$p(X=x_i)=\frac{c_i}{N}$$

• Fraction of the total number of points that fall in column *i*.

- 4 同下 4 三下 4 三下

Sum Rule

Number of instances in column *i* is just the sum of the number of instances in each cell of that column, we have $c_i = \sum_i n_{ij}$ and hence

$$p(X = x_i) = \frac{C_i}{N}$$
$$= \frac{\sum_{j=1}^{L} n_{ij}}{N}$$
$$= \sum_{j=1}^{L} p(X = x_i, Y = y_j).$$

Sometimes also called *marginal* probability, because it is obtained by marginalizing, or summing out, the other variables (here Y).

イロト イポト イヨト イヨト

Product Rule

Now we are interested in finding the fraction of those points in column i that fall in cell i, j, that is

$$p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}.$$

By means of this *conditional* probability, we can derive the following relationship

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N} = \frac{n_{ij}}{c_i} \cdot \frac{c_i}{N}$$
$$= p(Y = y_j | X = x_i) p(X = x_i)$$

which is the product rule of probabilities.

A B F A B F

Illustrative Example



Course IN2207

SS2011 30 / 256

Sum and Product Rule (Summary)

Rules of Probability



- 4 同 6 4 日 6 4 日 6

Bayes' theorem

Using symmetry property p(X, Y) = p(Y, X) one can derive the following relation between conditional probabilities

$$p(X, Y) = p(Y, X)$$

$$p(Y|X)p(X) = p(X|Y)p(Y)$$

$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$

where the denominator can be expressed by means of the sum rule as

$$p(X) = \sum_{Y} p(X, Y) = \sum_{Y} p(Y, X)$$
$$= \sum_{Y} p(X|Y)p(Y)$$

T.Stibor (TUM)

SS2011 32 / 256

A B K A B K

Bishop's Bayes Example



The red box contains 6 oranges and 2 apples, the blue box contains 1 orange and 3 apples. Suppose we pick the red box 40 % of the time and the blue box 60 % of the time. Thanks to Bayes' theorem we can answer questions such as:

- What is the overall probability that we pick an apple?
- Given that we have chosen an orange, what is the probability that the box we chose was the blue one?

Bishop's Bayes Example (cont.)

For the sake of clarity let's introduce random variables B for box and F for fruit. B can take one of the two possibilities B = r (for red) and B = b (for blue); and F = o (for orange) and F = a (for apple).

The prior probability of selecting the red box is

$$p(B=r)=\frac{4}{10}$$

and of selecting the blue box

$$p(B=b)=\frac{6}{10}$$

Bishop's Bayes Example (cont.)

From given information we can write out all four *conditional probabilities* of given the selected box and picking the type of fruit.

$$p(F = a|B = r) = \frac{1}{4}$$
$$p(F = o|B = r) = \frac{3}{4}$$
$$p(F = a|B = b) = \frac{3}{4}$$
$$p(F = o|B = b) = \frac{1}{4}$$



Bishop's Bayes Example (cont.)

Back to our question: What is the overall probability that we pick an apple?

$$p(F = a) = p(F = a|B = r)p(B = r) + p(F = a|B = b)p(B = b)$$
$$= \frac{1}{4} \cdot \frac{4}{10} + \frac{3}{4} \cdot \frac{6}{10} = \frac{11}{20}$$

from this it follows that $p(F = o) = 1 - \frac{11}{20} = \frac{9}{20}$. Although there are more oranges in total, picking an apple is more likely.

Back to our second question: Given that we have chosen an orange, what is the probability that the box we chose was the blue one, that is p(B = b|F = o)?
Bishop's Bayes Example (cont.)

$$p(B = b|F = o) = \frac{p(F = o|B = b)p(B = b)}{p(F = o)} = \frac{1}{4} \cdot \frac{6}{10} \cdot \frac{20}{9} = \frac{1}{3}$$

and that we chose the red box:

$$p(B = r|F = o) = \frac{p(F = o|B = r)p(B = r)}{p(F = o)} = \frac{3}{4} \cdot \frac{4}{10} \cdot \frac{20}{9} = \frac{2}{3}$$

$$p(B|F) = rac{p(F|B)p(B)}{p(F)}, ext{ where } p(F) = \sum_{B} p(F|B)p(B)$$

T.Stibor (TUM)

Course IN2207

▲ 王 · つへで
 SS2011 37 / 256

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Bayes' theorem Bishop's Example

- Probability p(B) is the prior probability because it is the probability before we observe the identify of the fruit.
- After telling us the fruit is e.g. an orange, we can use Bayes' theorem to compute the probability p(B|F) which is called the *posterior probability* because it is the probability *after* we have observed *F*.

Bayes' theorem in words

posterior \propto likelihood imes prior

Independence of Random Variables

• If random variables X and Y are independent, then joint distribution factorizes into product of marginals, that is

$$p(X,Y)=p(X)p(Y).$$

• From product rule one can see

$$p(Y|X) = p(Y),$$

that is, the conditional distribution of Y given X is independent of the value of X.

Naive Bayes Classifier

Given feature variables F_1, F_2, \ldots, F_n and a class variable C. The Bayes' theorem states

$$p(C | F_1, F_2, ..., F_n) = \frac{p(C) p(F_1, F_2, ..., F_n | C)}{p(F_1, F_2, ..., F_n)}.$$

Assuming that each feature F_i is conditionally independent of every other feature F_j for $i \neq j$ one obtains

$$p(C | F_1, F_2, ..., F_n) = \frac{p(C) \prod_{i=1}^n p(F_i | C)}{p(F_1, F_2, ..., F_n)}.$$

The denominator serves as a scaling factor and can be omitted in the final classification rule

$$\underset{c}{\operatorname{argmax}} p(C=c) \prod_{i=1}^{n} p(F_i = f_i \mid C = c).$$

T.Stibor (TUM)

Naive Bayes Classifier Example

Dataset taken from Tom Mitchell's book: Machine Learning

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Course IN2207

< ロ > < 同 > < 回 > < 回 > < 回

Naive Bayes Classifier Example (cont.)

Should we *play tennis* when given the following information (features F_1, F_2, F_3, F_4)

$$(F_1 = \text{sunny}, F_2 = \text{cool}, F_3 = \text{high}, F_4 = \text{strong})$$

We want to compute

$$\underset{c}{\operatorname{argmax}} p(C=c) \prod_{i=1}^{n} p(F_i = f_i \mid C = c).$$

Let's do it
$$p(C = Yes)p(F_1 = sunny|C = Yes)p(F_2 = cool|C = Yes)$$

 $p(F_3 = high|C = Yes)p(F_4 = strong|C = Yes) = 0.005$

$$p(C = \text{No})p(F_1 = \text{sunny}|C = \text{No})p(F_2 = \text{cool}|C = \text{No})$$
$$p(F_3 = \text{high}|C = \text{No})p(F_4 = \text{strong}|C = \text{No}) = 0.021$$

Hence, we better play not tennis.

A B K A B K

Probability Densities

- We considered probabilities defined over discrete sets of events.
- Consider now probabilities with respect to continuous variables.

If the probability of a real-valued variable x falling in the interval $(x, x + \delta x)$ is given by $p(x)\delta x$ for $\delta x \to 0$, then p(x) is called the *probability density* over x.

• The *probability* that x will lie in an interval (a, b) is given by

$$p(x \in (a, b)) = \int_a^b p(x) dx.$$

Probability Densities (cont.)

Probability density p(x) must satisfy the two conditions

$$p(x) \geq 0$$

 $\int_{-\infty}^{\infty} p(x) dx = 1$



Probability density can be expressed as the derivative of a cumulative distribution function P(x).

Multivariate Probability Densities

- Having a vector x = (x₁, x₂,..., x_D) of continuous variables, we can define the joint probability density p(x).
- Multivariate probability density must satisfy

٠

$$p(\mathbf{x}) \geq 0 \ \int_{-\infty}^{\infty} p(\mathbf{x}) d\mathbf{x} = 1$$

• Sum and product rule, as well as Bayes' theorem apply equally. If x and y are two real variables, then

$$p(x) = \int p(x,y) dy$$

$$p(x,y) = p(y|x)p(x).$$

Expectations

Average value of some function f(x) under a probability distribution p(x) is called the *expectation* of f(x), denoted as 𝔼[f].

For discrete distribution it is

$$\mathbb{E}[f] = \sum_{x} p(x)f(x).$$

Average is weighted by the relative probabilities of the different values of *x*. Continues variables are expressed in terms of an integration w.r.t. probability density

$$\mathbb{E}[f] = \int p(x)f(x)dx.$$

・ 同 ト ・ ヨ ト ・ ヨ ト

Expectations (cont.)

Given a finite number N of points drawn from a probability distribution or probability density, the expectation can be approximated as

$$\mathbb{E}(f)\approx \frac{1}{N}\sum_{n=1}^{N}f(x_n).$$

Subscript indicates which variable is being averaged over

$$\mathbb{E}_{x}[f(x,y)] = \sum_{x} p(x,y)f(x,y).$$

Conditional expectation w.r.t. a conditional distribution

$$\mathbb{E}_{x}[f \mid y] = \sum_{x} p(x \mid y) f(x).$$

Properties of Expectations

- If c is a constant, then $\mathbb{E}[c] = c$.
- $\mathbb{E}[ax] = a\mathbb{E}[x].$
- $\mathbb{E}[x+y] = \mathbb{E}[x] + \mathbb{E}[y].$
- $\mathbb{E}[x+c] = \mathbb{E}[x] + c$.

(B)

Variance and Covariance

Variance of f(x) is defined as

$$\operatorname{var}[f] = \mathbb{E}\left[(f(x) - \mathbb{E}[f(x)])^2\right]$$

and provides a measure of how much variability there is in f(x) around its mean value $\mathbb{E}[f(x)]$. Variance can also be written in terms of the expectations of f(x) and $f(x)^2$

$$\operatorname{var}[f] = \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2.$$

The variance of variable x itself is given by

$$\operatorname{var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2.$$

Variance and Covariance (cont.)

For two random variables x and y, the covariance is defined as

$$\begin{aligned} \mathsf{cov}[x,y] &= & \mathbb{E}_{x,y}\left[\{x - \mathbb{E}[x]\}\{y - \mathbb{E}[y]\}\right] \\ &= & \mathbb{E}_{x,y}[xy] - \mathbb{E}[x]\mathbb{E}[y]. \end{aligned}$$

Expectation, Variance and Covariance Example



total 60

$$\mathbb{E}[x] = p(X=1)1 + p(X=2)2 + \ldots + p(X=9)9$$

= $\frac{3}{60}1 + \frac{6}{60}2 + \ldots + \frac{2}{60}9 = 4.95$

$$\mathbb{E}[y] = p(Y=1)1 + p(Y=2)2$$

= $\frac{34}{60}1 + \frac{26}{60}2 = 1.4333$

T.Stibor (TUM)

SS2011 51 / 256

3 🕨 🖌 3

Expectation, Variance and Covariance Example (cont.)

$$\operatorname{var}[x] = \mathbb{E}\left[(x - \mathbb{E}[x])^2\right]$$

= $\frac{3}{60}(1 - 4.95)^2 + \frac{6}{60}(2 - 4.95)^2 + \frac{8}{60}(3 - 4.95)^2 + \dots + (9 - 4.95)^2\frac{2}{60} = 4.5475$

$$var[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2$$

= $1^2 \frac{3}{60} + 2^2 \frac{6}{60} + 3^2 \frac{8}{60} + \dots + 8^2 \frac{6}{60} + 9^2 \frac{2}{60} - 4.95^2 = 4.5475$

T.Stibor (TUM)

Course IN2207

SS2011 52 / 256

< 3 > < 3

Binary Variables & Bernoulli Distribution

Consider single binary random variable $x \in \{0, 1\}$, e.g. x might describe outcome of flipping a coin, with x = 1 representing 'heads', and x = 0 representing 'tails'.

Probability of x = 1

$$p(x = 1 \mid \mu) = \mu$$
 where $0 \le \mu \le 1$.

Probability of x = 0

$$p(x=0\,|\,\mu)=1-\mu.$$

Probability distribution over x is

Bern
$$(x \mid \mu) = \mu^{x} (1 - \mu)^{1-x}$$
 (*Bernoulli* distribution).

Distribution is normalized and has mean $\mathbb{E}[x] = \mu$ and variance $var[x] = \mu (1 - \mu)$.

直 ト イ ヨ ト イ ヨ ト

Bernoulli Distribution & Maximum Likelihood

Given a dataset $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$ of observed values of x. Assume that observations are drawn independently from $p(x \mid \mu)$, so that

$$p(\mathcal{D} \mid \mu) = \prod_{n=1}^{N} p(x_n \mid \mu) = \prod_{n=1}^{N} \mu^{x_n} (1-\mu)^{1-x_n}$$

Suppose we observe dataset $\mathcal{D} = \{0, 1, 1, 0, 1, 0, 1, 1, 1, 0\}$ and want to maximize $p(\mathcal{D} | \mu)$, that is, find a parameter μ such that $p(\mathcal{D} | \mu)$ is maximal.



Observe that for $\mu = 0.6 = \frac{\sum_{n=1}^{N} 1\{x_n=1\}}{|\mathcal{D}|}$ the value of $p(\mathcal{D} \mid \mu)$ is maximal.

T.Stibor (TUM)

Course IN2207

SS2011 54 / 256

Bernoulli Distribution & Maximum Likelihood (cont.)

Analytically we can set the derivative of $\ln p(\mathcal{D} \mid \mu)$ with respect to μ to zero.

$$\ln p(\mathcal{D} \mid \mu) = \sum_{n=1}^{N} \ln p(x_n \mid \mu) = \sum_{n=1}^{N} \{x_n \ln \mu + (1 - x_n) \ln(1 - \mu)\}.$$

$$\frac{\partial \ln p(\mathcal{D} \mid \mu)}{\partial \mu} = 0 \quad \Rightarrow \mu_{\mathsf{ML}} = \frac{1}{N} \sum_{n=1}^{N} x_n.$$

The maximum likelihood estimator μ_{ML} is also known as *sample mean* and is the fraction of observations of heads in the data set D, that is,

$$\mu_{\mathsf{ML}} = \frac{\sum_{n=1}^{N} \mathbb{1}\{x_n = 1\}}{|\mathcal{D}|}.$$

T.Stibor (TUM)

SS2011 55 / 256

Binomial Distribution

Suppose we toss a coin N times and are interested in the probability of obtaining m heads.

$$\mathsf{Bin}(m \mid \mathsf{N}, \mu) = \left(egin{array}{c} \mathsf{N} \ m \end{array}
ight) \mu^m (1-\mu)^{\mathsf{N}-m}.$$

For each observation the mean and variance are given by

$$\mathbb{E}[m] \equiv \sum_{m=0}^{N} m \operatorname{Bin}(m \,|\, N, \mu) = N \,\mu \quad (\text{sum of means})$$

$$\operatorname{var}[m] \equiv \sum_{m=0}^{N} (m - \mathbb{E}[m])^2 \operatorname{Bin}(m \mid N, \mu) = N \mu (1 - \mu) \quad (\text{sum of variances}).$$

Binomial Distribution Example

Suppose N = 3, $\mu = 1/2$ and m = 2, that is, we toss a fair coin three times and are interested in the probability of obtaining two times head.

The event space is $\Omega = \{H, T\}^3$ and the events of interests are $\mathcal{E} = \{\{H, H, T\}, \{H, T, H\}, \{T, H, H\}\}$. Observe that

$$\left(\begin{array}{c}3\\2\end{array}\right) = |\mathcal{E}|$$

and the probability of tossing two heads in three trials

$$3(1/2)^2(1-1/2)^1 = \begin{pmatrix} 3\\2 \end{pmatrix} (1/2)^2(1-1/2)^1.$$

T.Stibor (TUM)

ヘロト 人間 ト 人 ヨ ト 一

Multinomial Distribution

So far we have considered random variables that can take two possible values (e.g. head or tail). The multinomial distribution encounters discrete variables that can take on one of K possible mutually exclusive states.

Denote *N* the total number of observations, μ parameter vector with $0 \le \mu_k \le 1$ and $\sum_{k=1}^{K} \mu_k = 1$. The *multinomial distribution* has the form

$$\mathsf{Mult}(m_1, m_2, \ldots, m_K \,|\, oldsymbol{\mu}, N) = \left(egin{array}{c} N \ m_1 m_2 \ldots m_K \end{array}
ight) \prod_{k=1}^K \mu_k^{m_k},$$

where m_k denotes the number of times outcome number k was observed over the N observations (trials) and

$$\left(\begin{array}{c}N\\m_1m_2\ldots m_K\end{array}\right)=\frac{N!}{m_1!m_2!\ldots m_K!}$$

(number of ways of partitioning N objects into K groups of size m_1, m_2, \ldots, m_K).

T.Stibor (TUM)

SS2011 58 / 256

. .

Multinomial Distribution (Example)

Given an urn with 10 marbles, 2 red, 3 green and 5 blue (K = 3). We select randomly N = 4 from the urn with replacement. What is the probability of selecting $m_1 = 0$ red, $m_2 = 2$ green and $m_3 = 2$ blue marbles.

$$p(red) = 2/10, p(green) = 3/10, p(blue) = 5/10.$$

 $Mult(m_1 = 0, m_2 = 2, m_3 = 2 | \mu_1 = 2/10, \mu_2 = 3/10, \mu_3 = 5/10, N = 4)$

$$=\frac{4!}{0!\,2!\,2!}(0.2)^0\,(0.3)^2\,(0.5)^2=0.135,$$

where $m_1 + m_2 + m_3 = N = 4$.

イロト イポト イヨト イヨト

Gaussian Distribution

For a single real-valued variable x, the Gaussian distribution is defined by

$$\mathcal{N}(x|\mu,\sigma) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\}$$

where μ is the mean and σ^2 the variance (square root of σ^2 is called standard deviation).

Gaussian distribution satisfies

$$\mathcal{N}(x|\mu,\sigma) \geq 0$$
 $\int_{-\infty}^{\infty} \mathcal{N}(x|\mu,\sigma) \, dx = 1$

and therefore satisfies the two requirements for a valid probability density.

Gaussian Distribution (cont.)



SS2011 61 / 256

Multivariate Gaussian Distribution

Gaussian distribution defined over D-dimensional vector \mathbf{x} of continuous variables is given by

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}$$
$$\boldsymbol{\mu} = \begin{bmatrix} 0\\0 \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0\\0 & 1 \end{bmatrix}$$



- μ is a *D*-dim. mean vector.
- Σ is a D × D covariance matrix, |Σ| denotes the determinant of Σ.

Gaussian Distribution & Maximum Likelihood

- Given data set $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$ of observations.
- Suppose that observations are drawn independently from a Gaussian distribution whose mean μ and variance σ^2 are *unknown*.
- Goal: infer these parameters from the data set.

Data ${\cal X}$ is independent and identically distributed, hence likelihood of the data, given μ and σ^2 is

$$p(\mathcal{X} \mid \mu, \sigma^2) = \prod_{n=1}^N \mathcal{N}(x_n \mid \mu, \sigma^2).$$

Log likelihood of $p(\mathcal{X} | \mu, \sigma^2)$

$$\ln p(\mathcal{X} \mid \mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^{N} (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi)$$

Gaussian Distribution & Maximum Likelihood (cont.) Maximizing $\ln p(\mathcal{X} | \mu, \sigma^2)$, that is

$$rac{\partial}{\partial \mu} \ln p(\mathcal{X} \,|\, \mu, \sigma^2) = 0$$

gives solution

$$\mu_{\mathsf{ML}} = rac{1}{N} \sum_{n=1}^{N} x_n$$
 sample mean

and analog

$$\frac{\partial}{\partial \sigma} \ln p(\mathcal{X} \mid \mu, \sigma^2) = 0$$

gives sample variance

$$\sigma_{\rm ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{\rm ML})^2$$

Properties of Gaussian Distribution

Expectation of $\mathbf{x} \in \mathbb{R}^D$ is

$$\mathbb{E}[\mathbf{x}] = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\mathbf{\Sigma}|^{1/2}} \int \exp\left\{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \mathbf{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\right\} \mathbf{x} \, d\mathbf{x} \\ = \boldsymbol{\mu}$$

and covariance

$$cov[\mathbf{x}] = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T]$$

$$cov[\mathbf{x}] = \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T]$$

$$= \boldsymbol{\Sigma}.$$

Covariance matrix Σ is always symmetric and positive semidefinite. Elements σ_{ii} are variances of x_i , elements σ_{ij} are covariances of x_i and x_j . If x_i and x_j are statistically independent, then $\sigma_{ij} = 0$.

Two-dim. Gaussian Distribution

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{D/2} |\mathbf{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}$$
$$\boldsymbol{\mu} = \begin{bmatrix} 0\\0 \end{bmatrix} \quad \mathbf{\Sigma} = \begin{bmatrix} 1 & 0\\0 & 1 \end{bmatrix}$$



T.Stibor (TUM)

SS2011 66 / 256

Two-dim. Gaussian Distribution



T.Stibor (TUM)

Course IN2207

SS2011 67 / 256

Two-dim. Gaussian Distribution (cont.)



T.Stibor (TUM)

Course IN2207

SS2011 68 / 256

Covariance Matrix

Covariance matrix Σ is symmetric and positive definite. Example:

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \mathbf{\Sigma} = \mathbf{\Sigma}^{T}.$$
$$x, y \in \mathbb{R}, \begin{bmatrix} x & y \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} x & y \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = x^{2} + y^{2} \ge 0, \text{ for all } x, y \in \mathbb{R}$$





(3)

Covariance Matrix

Eigenvector and eigenvalue of covariance matrix $\mathbf{\Sigma} \in \mathbb{R}^{D imes D}$

$$\mathbf{\Sigma}\mathbf{u}_d = \lambda_d \mathbf{u}_d, \quad d = 1, 2, \dots, D.$$

Covariance matrix can be expressed as an expansion in terms of its eigenvectors in the form

$$\mathbf{\Sigma} = \sum_{d=1}^{D} \lambda_d \mathbf{u}_d \mathbf{u}_d^{\mathsf{T}}$$

and the inverse matrix

$$\mathbf{\Sigma}^{-1} = \sum_{d=1}^{D} \frac{1}{\lambda_d} \mathbf{u}_d \mathbf{u}_d^T.$$

A B + A B +

Covariance Matrix and Eigenvectors (Example) $\boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \qquad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & 1/2 \\ 1/2 & 1 \end{bmatrix}$



T.Stibor (TUM)

Course IN2207

SS2011 71 / 256

Linear Combinations of Gaussian Distributions

Suppose $p(\mathbf{x}) \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, and **A** is *d*-by-*k* matrix and $\mathbf{y} = \mathbf{A}^T \mathbf{x}$ is a *k*-component vector, then

 $p(\mathbf{y}) \sim N(\mathbf{A}^{\mathsf{T}} \boldsymbol{\mu}, \mathbf{A}^{\mathsf{T}} \boldsymbol{\Sigma} \mathbf{A}).$

When k = 1 and **A** is a unit-length vector **a**, $y = \mathbf{a}^T \mathbf{x}$ is a scalar that represents the projection of **x** onto a line in the direction of **a**, where $\mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a}$ is the variance of the projection of **x** onto **a**.



Course IN2207
Mahalanobis Distance

Quantity

$$r^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

is called the squared *Mahalanobis Distance* from **x** to μ . If $\Sigma = I$, then r^2 is the squared Euclidean distance. Any points lying on the same contour line have equal distance to μ .



$$[x_1 x_2] \begin{bmatrix} 1 & 1/2 \\ 1/2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = x_1(x_1 + 1/2 x_2) + x_2(1/2 x_1 + x_2) = x_1^2 + x_1 x_2 + x_2^2.$$

 x_1

Course IN2207

Bayesian Decision Theory

- How to make an optimal decision given the appropriate probabilities?
- Minimize the error of assigning **x** to the wrong class.
- Intuitively we would choose the class having the higher posterior probability.

An error occurs when x belonging to class C_1 is assigned to class C_2 or vice versa. The probability of this occuring is given by:

$$p(\text{error}) = p(\mathbf{x} \in \mathcal{R}_1, \mathcal{C}_2) + p(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1)$$
$$= \int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) d\mathbf{x}$$

Error Probabilities in Bay. Decision Theory



Course IN2207

Error Probabilities in Bay. Decision Theory



Course IN2207

▲ 王 シ へ ペ
 SS2011 76 / 256

• • = • • = •

Error Probabilities in Bay. Decision Theory



Course IN2207

▲ 王 シ へ ペ
 SS2011 77 / 256

伺下 イヨト イヨト

Bayesian Estimation (Class-Conditional Densities)

In Bayesian parameter estimation, the parameter is not fixed as in ML, but is a random variable. The training data is used to convert a distribution on this variable into a posterior density.

Let D denote a training sample, our goal is to compute the posterior p(C_i | x, D).

Bayes formula:

$$p(C_i \mid x, \mathcal{D}) = \frac{p(x \mid C_i, \mathcal{D})p(C_i \mid \mathcal{D})}{\sum_{i=1}^{I} p(x \mid C_i, \mathcal{D})p(C_i \mid \mathcal{D})}.$$

Training sample ${\mathcal D}$ can be used to help us determining the class-conditional and prior density.

To simplify, we assume that prior probabilities are known or obtainable from some calculations, thus we substitute $P(C_i) = P(C_i | D)$.

ヘロア 人間ア 人間ア 人間ア

Consider supervised case: Separate training data by class into c subsets $\mathcal{D}_1, \mathcal{D}_2, \ldots, \mathcal{D}_I$ with sample in \mathcal{D}_i belonging to class C_i . This assumption allows us to write the previous equation as

$$p(C_i | x, \mathcal{D}) = \frac{p(x | C_i, \mathcal{D})p(C_i)}{\sum_{i=1}^{I} p(x | C_i, \mathcal{D})p(C_i)}.$$

We have now *I* separate problems of the following form: Use a set of \mathcal{D} of samples drawn independently according to a fixed but unknown probability distribution p(x) to determine $p(x|\mathcal{D})$.

Bayesian Parameter Estimation

- Desired density p(x) is unknown, we however assume that it has a known parametric form. Unknown is just the value of the parameter vector Θ. That is p(x | Θ) is known.
- Any information we have about Θ prior to observing the sample is contained in prior density $p(\Theta)$.

Observation of the sample \mathcal{D} converts this to a posterior density $p(\Theta | \mathcal{D})$, which is close (sharply peaked) about the true value of Θ .

Note that we have converted our problem if learning a probability density function to one of estimating a parameter vector.

Recall that our goal is to compute p(x | D), which is as close as we can come to obtain the unknown density p(x).

イロト 不得下 イヨト イヨト 二日

Bayesian Parameter Estimation (cont.)

We obtain p(x | D) by integrating the joint density $p(x, \Theta | D)$ over Θ :

$$p(x \mid \mathcal{D}) = \int p(x, \Theta \mid \mathcal{D}) d\Theta,$$

= $\int p(x \mid \Theta, \mathcal{D}) p(\Theta \mid \mathcal{D}) d\Theta$

If $p(\boldsymbol{\Theta} \mid \mathcal{D})$ peaks very sharply about some value $\hat{\boldsymbol{\Theta}}$, we obtain $p(x \mid \mathcal{D}) \approx p(x \mid \boldsymbol{\Theta})$.

< 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Bayesian Parameter Estimation (μ)

Calculate the posterior density $p(\boldsymbol{\Theta}|\mathcal{D})$ and the desired density $p(x|\mathcal{D})$ for the case where $p(x|\mu) \sim N(\mu, \sigma^2)$, where the only unknown quantity is $\boldsymbol{\Theta} \equiv \mu$.

Whatever prior knowledge we might have about μ can be expressed by a known prior density $p(\mu)$, where me make the assumption that

 $p(\mu) \sim N(\mu_0, \sigma_0^2)$, where both μ_0 and σ_0 are known.

Intuitively, μ_0 represents our best prior guess for μ , and σ_0^2 measures our uncertainty about this guess.

イロト 不得下 イヨト イヨト 二日

Bayesian Parameter Estimation (μ cont.)

Let
$$\mathcal{D} = \{x_1, x_2, \dots, x_N\}.$$

$$p(\mu|\mathcal{D}) = \frac{p(\mathcal{D}|\mu)p(\mu)}{\int p(\mathcal{D}|\mu)p(\mu) d\mu} = \frac{p(\mathcal{D}|\mu)p(\mu)}{p(\mathcal{D})}$$

$$= \alpha \prod_{n=1}^{N} p(x_n|\mu)p(\mu),$$

where α is a normalization factor that depends on \mathcal{D} but is independent of μ . Because $p(x_k|\mu) \sim N(\mu, \sigma^2)$ and $p(\mu) \sim N(\mu_0, \sigma_0^2)$ we obtain

$$p(\mu|\mathcal{D}) = \alpha \prod_{n=1}^{N} \underbrace{\frac{1}{\sqrt{2\pi\sigma}} \exp\left[-\frac{1}{2}\left(\frac{x_n - \mu}{\sigma}\right)^2\right]}_{p(\mu)} \underbrace{\frac{1}{\sqrt{2\pi\sigma_0}} \exp\left[-\frac{1}{2}\left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right]}_{p(\mu)}$$

T.Stibor (TUM)

SS2011 83 / 256

イヨト イヨト

Bayesian Parameter Estimation (μ cont.)

Additional manipulations give

$$p(\mu|\mathcal{D}) = \alpha' \exp\left[-\frac{1}{2}\left(\sum_{n=1}^{N}\left(\frac{\mu-x_n}{\sigma}\right)^2 + \left(\frac{\mu-\mu_0}{\sigma_0}\right)^2\right)\right]$$
$$= \alpha'' \exp\left[-\frac{1}{2}\left[\left(\frac{N}{\sigma^2} + \frac{1}{\sigma_0^2}\right)\mu^2 - 2\left(\frac{1}{\sigma^2}\sum_{n=1}^{N}x_n + \frac{\mu_0}{\sigma_0^2}\right)\mu\right]\right] (1)$$

If we write $p(\mu|D) \sim N(\mu_N, \sigma_N^2)$, then μ_N and σ_N^2 can be found by equating coefficients in (1) with corresponding coefficients in the generic Gaussian of the form

$$p(\mu|\mathcal{D}) = rac{1}{\sqrt{2\pi\sigma_N}} \exp\left[-rac{1}{2}\left(rac{\mu-\mu_N}{\sigma_N}
ight)^2
ight]$$

T.Stibor (TUM)

SS2011 84 / 256

- 4 同下 4 三下 4 三下

Bayesian Parameter Estimation (μ cont.)

Solving for μ_N and σ_N^2 gives

$$\mu_N = \left(\frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2}\right)\hat{\mu_N} + \frac{\sigma^2}{N\sigma_0^2 + \sigma^2}\mu_0, \text{ where } \hat{\mu_N} = \frac{1}{N}\sum_{n=1}^N x_n$$

and

$$\sigma_N^2 = \frac{\sigma_0^2 \sigma^2}{N \sigma_0^2 + \sigma^2}.$$

Intuitively, μ_N represents our best guess for μ after observing a sample of size N, and σ_N^2 measures the uncertainty about this guess. Note that σ_N^2 decreases monotonically with $N \to \infty$, that is, each additional observation decreases our uncertainty about the true nature of μ . As N increases, $p(x|\mathcal{D})$ becomes more and more sharply peaked.

Bayesian Parameter Estimation (Example)

We want to apply Bayesian parameter estimation (inference) for parameter μ when given a sample $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$ where $x_n \sim N(0.8, 0.1)$ and N = 10. The prior is $p(\mu) \sim N(0, 1)$.



Ν	=	1
Ν	=	2
Ν		3
Ν	=	4
Ν	=	10

T.Stibor (TUM)

Course IN2207

S2011 86 / 256

History of Neural Networks



Note, this historical overview is far from being complete (see books for detailed historical overview)

T.Stibor (TUM)

Course IN2207

SS2011 87 / 256

Method of Steepest Descent

Let $E(\mathbf{w})$ be a continuously differentiable function of some unknown (weight) vector \mathbf{w} .

Find an optimal solution \boldsymbol{w}^{\star} that satisfies the condition

$$E(\mathbf{w}^{\star}) \leq E(\mathbf{w}).$$

The necessary condition for optimality is

$$\nabla E(\mathbf{w}^{\star}) = \mathbf{0}.$$

Let us consider the following *iterative* descent:

Start with an initial guess $\mathbf{w}^{(0)}$ and generate sequence of weight vectors $\mathbf{w}^{(1)}, \mathbf{w}^{(2)}, \ldots$ such that

$$E(\mathbf{w}^{(i+1)}) \leq E(\mathbf{w}^{(i)}).$$

Steepest Descent Algorithm

$$\mathbf{w}^{(i+1)} = \mathbf{w}^{(i)} - \eta \nabla E(\mathbf{w}^{(i)})$$

where η is a positive constant called learning rate.

At each iteration step the algorithm applies the correction

$$\Delta \mathbf{w}^{(i)} = \mathbf{w}^{(i+1)} - \mathbf{w}^{(i)}$$
$$= -\eta \nabla E(\mathbf{w}^{(i)})$$

Steepest descent algorithm satisfies:

$$E(\mathbf{w}^{(i+1)}) \leq E(\mathbf{w}^{(i)}),$$

to see this, use first-order Taylor expansion around $\mathbf{w}^{(i)}$ to approximate $E(\mathbf{w}^{(i+1)})$ as $E(\mathbf{w}^{(i)}) + (\nabla E(\mathbf{w}^{(i)}))^T \Delta \mathbf{w}^{(i)}$.

Steepest Descent Algorithm (cont.)

$$E(\mathbf{w}^{(i+1)}) \approx E(\mathbf{w}^{(i)}) + (\nabla E(\mathbf{w}^{(i)}))^T \Delta \mathbf{w}^{(i)}$$

= $E(\mathbf{w}^{(i)}) - \eta \|\nabla E(\mathbf{w}^{(i)})\|^2$

For positive learning rate η , $E(\mathbf{w}^{(i)})$ decreases in each iteration step (for small enough learning rates).



T.Stibor (TUM)

Course IN2207

SS2011 90 / 256

Steepest Descent Algorithm Example

Black points denote different starting values. Learning rate η is properly chosen, however for starting value (1, 1), algorithm converges not to the global minimum. It follows steepest descent in the "wrong direction", in other words, gradient based algorithms are local search algorithms.



T.Stibor (TUM)

Course IN2207

SS2011 91 / 256

Steepest Descent Algorithm Example (cont.)

Learning rate $\eta = 1.0$ is too large, algorithm oscillates in a "zig-zag" manner or "overleap" the global minimum.



ヨト イヨト

Steepest Descent Algorithm Example (cont.)

Learning rate $\eta = 0.005$ is too small, algorithm converges "very slowly".



T.Stibor (TUM)

Course IN2207

SS2011 93 / 256

Single-Layer Network



Equivalent notation:

$$y(\mathbf{x}) = \widetilde{\mathbf{w}}^T \widetilde{\mathbf{x}} = \sum_{i=0}^a w_i x_i$$

٦

where $\widetilde{\mathbf{w}} = (w_0, \mathbf{w})$ and $\widetilde{\mathbf{x}} = (1, \mathbf{x})$.

SS2011 94 / 256

3 1 4 3

Linear Classifier

• Linear classifiers are single layer neural networks.



SS2011 95 / 256

< 同 > < 三 > < 三 > <

Linear Classifier & Dot Product



- What about the vector
 w = (w₁, w₂) = (-2, 1)?
- Vector **w** is perpendicular to the line $-2x_1 + 1x_2 = 0$.
- Let us calculate the dot product of **w** and **x**.

The dot product is defined as $w_1x_1 + w_2x_2 + \ldots + w_dx_d \stackrel{\text{def}}{=} \langle \mathbf{w}, \mathbf{x} \rangle$, for some $d \in \mathbb{N}$.

In our example d = 2 and we obtain $-2 \cdot 1 + 1 \cdot 2 = 0$.

A B K A B K

Linear Classifier & Dot Product (cont.)

Let us consider the *weight* vector $\mathbf{w} = (3, 0)$ and vector $\mathbf{x} = (2, 2)$.



Geometric interpretation of the dot product: Length of the projection of **x** onto the unit vector $\mathbf{w}/||\mathbf{w}||$.

T.Stibor (TUM)

SS2011 97 / 256

Linear Classifier & Two Half-Spaces



The x-space is separated in two half-spaces.

T.Stibor (TUM)

SS2011 98 / 256

Linear Classifier & Dot Product (cont.)

- Observe, that $w_1x_1 + w_2x_2 = 0$ implies, that the separating line always goes through the origin.
- By adding an offset (bias), that is $w_0 + w_1x_1 + w_2x_2 = 0 \Leftrightarrow x_2 = -\frac{w_1}{w_2}x_1 - \frac{w_0}{w_2} \equiv y = mx + b$, one can shift the line arbitrary.



Linear Classifier & Single Layer NN



Given data which we want to separate, that is, a sample $\mathcal{X} = \{(\mathbf{x}_1, t_1), (\mathbf{x}_2, t_2), \dots, (\mathbf{x}_N, t_N)\} \in \mathbb{R}^d \times \{-1, +1\}.$

How to determine the proper values of **w** such that the "minus" and "plus" points are separated by $y(\mathbf{x})$? Infer the values of **w** from the data by some learning algorithm.

Perceptron

Note, so far we have not seen a method for finding the weight vector \mathbf{w} to obtain a linearly separation of the training set.

Let g(a) be (sign) activation function

$$g(a) = \left\{ egin{array}{cc} -1 & ext{if } a < 0 \ +1 & ext{if } a \geq 0 \end{array}
ight.$$

and decision function

$$g(\langle \mathbf{w}, \mathbf{x} \rangle) = g\left(\sum_{i=0}^d w_i x_i\right).$$

Note: x_0 is set to +1, that is, $\mathbf{x} = (1, x_1, \dots, x_d)$. Training pattern consists of $(\mathbf{x}, t) \in \mathbb{R}^{d+1} \times \{-1, +1\}$

(4 間) トイヨト イヨト

Perceptron Learning Algorithm

```
input : (\mathbf{x}_1, t_1), \ldots, (\mathbf{x}_N, t_N) \in \mathbb{R}^{d+1} \times \{-1, +1\}, \eta \in \mathbb{R}_+, \text{max.epoch} \in \mathbb{N}
output: w
begin
    Randomly initialize w
    epoch \leftarrow 0
    repeat
         for i \leftarrow 1 to N do
          epoch \leftarrow epoch + 1
    until (epoch = max.epoch) or (no change in \mathbf{w})
    return w
end
```

.

Training the Perceptron (cont.)

Geometrical explanation: If **x** belongs to {+1} and $\langle \mathbf{w}, \mathbf{x} \rangle < 0 \Rightarrow$ angle between **x** and **w** is greater than 90°, rotate **w** in direction of **x** to bring missclassified **x** into the positive half space defined by **w**. Same idea if **x** belongs to {-1} and $\langle \mathbf{w}, \mathbf{x} \rangle \ge 0$.



	·
I Stubor /	
1.5000 (

Course IN2207

SS2011 103 / 256

.

Perceptron Error Reduction

Recall: missclassifcation results in:

$$\mathbf{w}^{\text{new}} = \mathbf{w} + \eta \mathbf{x} t,$$

this reduces the error since

$$-\mathbf{w}^{\text{new}}(\mathbf{x} t)^{T} = -\mathbf{w}(\mathbf{x} t)^{T} - \underbrace{\eta}_{>0} \underbrace{(\mathbf{x} t)(\mathbf{x} t)^{T}}_{\|\mathbf{x}t\|^{2} > 0}$$
$$< -\mathbf{w}^{T} \mathbf{x} t$$

How often one has to cycle through the patterns in the training set?

• A finite number of steps?

Perceptron Convergence Theorem

Proposition

Given a finite and linearly separable training set. The perceptron converges after some finite steps [Rosenblatt, 1962].

Perceptron Algorithm (R-code)

```
*******************
perceptron <- function(w,X,t,eta,max.epoch) {</pre>
N \leq nrow(X)/2;
 epoch <- 0;
 repeat {
   w.old <- w:
   for (i in 1:(2*N)) {
     if (t[i]*y(X[i,],w) <= 0)
       w \leftarrow w + eta * t[i] * X[i.]:
   }
   epoch <- epoch + 1;</pre>
   if ( identical(w.old,w) || epoch = max.epoch ) {
     break; # terminate if no change in weights or max.epoch
   }
 3
 return (w);
                                      ▲圖▶ ▲屋▶ ▲屋▶
                                                  = 200
    T.Stibor (TUM)
                                             SS2011
                                                   106 / 256
```

Perceptron Algorithm Visualization



One epoch

terminate if no change in ${\boldsymbol{w}}$

Course IN2207

Perceptron Algorithm Visualization



One epoch

terminate if no change in ${\boldsymbol{w}}$

Course IN2207
Least Mean Square

Let us consider the weight correction in terms of an error function $E^{(i)} = \frac{1}{2} (\underbrace{y^{(i)}}_{g(\mathbf{w}^T \mathbf{x})} - t^{(i)})^2$, where $g(\cdot)$ is a differentiable function. Apply gradient descent rule

$$\mathbf{w}^{(i+1)} = \mathbf{w}^{(i)} - \eta \frac{\partial E^{(i)}}{\partial \mathbf{w}}, \text{ where } \frac{\partial E^{(i)}}{\partial \mathbf{w}} = \underbrace{(y^{(i)} - t^{(i)})}_{\delta^{(i)}} \mathbf{x}^{(i)}$$

gives change in weights

$$\Delta \mathbf{w} = -\eta \delta^{(i)} \mathbf{x}^{(i)} = -\eta \frac{\partial E^{(i)}}{\partial \mathbf{w}}$$

Delta rule \equiv {Adaline rule, Widrow-Hoff rule, Least Mean Square (LMS) }

Least Mean Square

Note, if we choose g(a) = a to be the linear activation function $\mathbf{w}^T \mathbf{x}$, then there exists a closed analytical solution (pseudo-inverse solution).

Let g(a) be a differentiable non-linear activation function, where $a = \mathbf{w}^T \mathbf{x}$.

$$rac{\partial E^{(i)}}{\partial \mathbf{w}} = \delta^{(i)} \mathbf{x}^{(i)}, ext{ where } \delta^{(i)} = g'(a)(y^{(i)} - t^{(i)})$$

gives change in weights

$$\Delta \mathbf{w} = -\eta \delta^{(i)} \mathbf{x}^{(i)} = -\eta \frac{\partial E^{(i)}}{\partial \mathbf{w}}$$

LMS Online/Batch Learning

Online learning:

• Update weight
$$\mathbf{w}^{(i+1)} = \mathbf{w}^{(i)} - \eta \frac{\partial E^{(i)}}{\partial \mathbf{w}}$$
 (pattern by pattern).

This type of online learning is also called *stochastic gradient descent*, it is an approximation of the true gradient.

Batch learning:

• Update weight $\mathbf{w}^{(i+1)} = \mathbf{w}^{(i)} - \eta \sum_{i=1}^{N} \frac{\partial E^{(i)}}{\partial \mathbf{w}}$ by computing derivatives for each pattern separately and then sum over all patterns.

Minimum Squared Error and Pseudoinverse

Recall that we want to minimize the squared error

. .

$$E(\mathbf{w}) = \sum_{i=1}^{N} \frac{1}{2} \left(y^{(i)} - t^{(i)} \right)^2$$
 where $y^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)}$

Let **X** be the $N \times \tilde{d}$ matrix where $\tilde{d} = d + 1$ and *i*th row denotes training pattern $\mathbf{x}^{(i) T}$, **w** is weight vector, **t** class label vector.

$$\begin{pmatrix} x_{10} & x_{11} & \cdots & x_{1d} \\ x_{20} & x_{21} & \cdots & x_{2d} \\ \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \\ x_{N0} & x_{N1} & \cdots & x_{Nd} \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{pmatrix} = \begin{pmatrix} t_0 \\ t_1 \\ \vdots \\ \vdots \\ t_N \end{pmatrix} \qquad \mathbf{Xw} = \mathbf{t}$$

T.Stibor (TUM)

Course IN2207

SS2011 112 / 256

MSE and Pseudoinverse (cont.)

Problem: find weight vector \mathbf{w} , that is, solve $\mathbf{X}\mathbf{w} = \mathbf{t}$.

If **X** is non-singular solve $\mathbf{w} = \mathbf{X}^{-1}\mathbf{t}$, however, if **X** is rectangular (which is usually the case), then there are more equations than unknowns, that is, the equation system is overdetermined.

Let us search for \mathbf{w} that minimizes the error

 $\mathbf{e} = \mathbf{X}\mathbf{w} - \mathbf{t}$

one approach is to minimize the squared length of the error vector e

$$J(\widetilde{\mathbf{w}}) = \|\mathbf{X}\mathbf{w} - \mathbf{t}\|^2 = \sum_{i=1}^{N} \left(\mathbf{w}^T \mathbf{x}^{(i)} - t^{(i)}\right)^2$$

MSE and Pseudoinverse (cont.)

Forming the gradient

$$\nabla J = \sum_{i=1}^{N} 2\left(\mathbf{w}^{T}\mathbf{x}^{(i)} - t^{(i)}\right)\mathbf{x}^{(i)} = 2\mathbf{X}^{T}(\mathbf{X}\mathbf{w} - \mathbf{t})$$

and setting ∇J to zero gives $\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{t}$. Observe that $\mathbf{X}^T \mathbf{X}$ is a $\tilde{d} \times \tilde{d}$ matrix which often is non-singular. In the non-singular case, one can solve \mathbf{w} uniquely as

$$\mathbf{w} = \left(\mathbf{X}^{\mathsf{T}} \mathbf{X} \right)^{-1} \mathbf{X}^{\mathsf{T}} \mathbf{t}$$
$$= \mathbf{X}^{\dagger} \mathbf{t}$$

The $\tilde{d} \times N$ matrix $\mathbf{X}^{\dagger} \equiv (\mathbf{X}^{T}\mathbf{X})^{-1}\mathbf{X}^{T}$ is called *pseudoinverse* of \mathbf{X} .

Linear Separability

Decision boundaries of single-layer networks are linear (hyperplanar in higher dimensions).

- Very restricted class of decision boundaries
- Examples:



XOR-Problem

Points are not linearly separable

Course IN2207

Probability for Linear Separability

- Probability that a random set of points will be linearly separable
- Suppose we have *N* points distributed at random in *d* dimensions in general position (not collinear)
- Randomly assign each of the points to one of the two classes C_1 and C_2 (with eq. probability)
- For N data points there are 2^N possible class assignments (dichotomies ≡ binary partitions)

Question: What fraction F(N, d) of these dichotomies is linearly separable?

Probability for Linear Separability (cont.)

$$F(N,d) = \begin{cases} 1 & \text{when } N \leq d+1 \\ \frac{1}{2^{N-1}} \sum_{i=0}^{d} \binom{N-1}{i} & \text{when } N \geq d+1 \end{cases}$$



If number of points is $\leq d + 1$, then any labeling leads to a separable problem.

T.Stibor (TUM)

Course IN2207

SS2011 117 / 256

Activation functions

Discrimination functions of the form $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ are simple linear functions of the input variables \mathbf{x} , where distances are measured by means of the dot product.

Let us consider the non-linear *logistic sigmoid* activation function $g(\cdot)$ for limiting the output to (0,1), that is,



Single-layer network with a logistic sigmoid activation function can also output posterior probabilities (rather than geometric distances).

イロト イポト イヨト イヨト

Activation functions (cont.)

Heaviside step function:

$$g(a) = \begin{cases} 0 & \text{if } a < 0 \\ 1 & \text{if } a \ge 0 \end{cases}$$



Hyperbolic tangent function:

$$g(a) = anh(a) = rac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)}$$

Note, $tanh(a) \in (-1,1)$



Course IN2207

1

Multi-Layer vs. Single-Layer Networks

Single-layer networks

- based on a linear combination of the input variables which is transformed by linear/non-linear activation function
- are limited in terms of the range of functions they can represent

Multi-layer networks

- consist of multiple layers and are capable of approximating any continuous functional mapping
- are compared to single-layer networks not so straightforward to train

Multi-Layer Network



Connection in first layer from input unit *i* to hidden unit *j* is denoted as w_{ji} . Connection from hidden unit *j* to output unit *k* is denoted as v_{ki} .

Multi-Layer Network (cont.)

Hidden unit *j* receives input

$$a_j = \sum_{i=1}^d w_{ji} x_i + w_{j0} = \sum_{i=0}^d w_{ji} x_i$$

and produces output

$$z_j = g(a_j) = g\left(\sum_{i=0}^d w_{ji}x_i\right).$$

Output unit k thus receives

$$a_k = \sum_{j=1}^M v_{kj} z_j + v_{k0} = \sum_{j=0}^M v_{kj} z_j$$

T.Stibor (TUM)

Multi-Layer Network (cont.)

and produces the final output

$$y_k = g(a_k) = g\left(\sum_{j=0}^M v_{kj} z_j\right) = g\left(\sum_{j=0}^M v_{kj} g\left(\sum_{i=0}^d w_{ji} x_i\right)\right)$$

Note that the activation function $g(\cdot)$ in the first layer can be different from those in the second layer (or other layers).

3 ∃ ≥

Multi-Layer Networks Example



Note: sometimes the layers of units are counted (here three layers), rather the layers of adaptive weights. In this course L-layer network is referred to a network with L layers of adaptive weights.

Course IN2207

SS2011 124 / 256

LMS Learning Rule for Multi-Layer Networks

- We have seen that the LMS learning rule is based on the gradient descent algorithm.
- The LMS learning rule worked because the error is proportional to the square difference between actual output y and target output t and can be evaluated for each output unit.
- In a multi-layer network we can use LMS learning rule on the hidden-to-output layer weights because target outputs are known.

Problem: we cannot compute the target outputs of the input-to-hidden weights because these values are unknown, or, to put it the other way around, how to update the weights in the first layer?

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Backpropagation (Hidden-to-Output Layer)

Recall that we want to minimize the error on training patterns between actual output y_k and target output t_k :

$$E = rac{1}{2} \sum_{k=1}^{K} (y_k - t_k)^2.$$

Backpropagation learning rule is based on gradient descent:

$$\Delta \mathbf{w} = -\eta \frac{\partial E}{\partial \mathbf{w}}$$
, component form $\Delta w_{st} = -\eta \frac{\partial E}{\partial w_{st}}$

Apply chain rule for differentiation:

$$\frac{\partial E}{\partial v_{kj}} = \frac{\partial E}{\partial a_k} \frac{\partial a_k}{\partial v_{kj}}$$

T.Stibor (TUM)

Course IN2207

- 4 同 6 4 日 6 4 日 6

Backprop. (Hidden-to-Output Layer) (cont.)

Gradient descent rule gives:

$$\begin{aligned} \Delta v_{kj} &= -\eta \frac{\partial E}{\partial v_{kj}} = -\eta (y_k - t_k) g'(a_k) z_j \\ &= -\eta \delta_k z_j \end{aligned}$$

where

$$\delta_k = (y_k - t_k)g'(a_k).$$

Observe that this result is identical to that obtained for LMS.

T Callen /	
I.SLIDOR (

Backpropagation (Input-to-Hidden Layer)

For the input-to-hidden connection we must differentiate with respect to the w_{ii} 's which are deeply embedded in

$$E = \frac{1}{2} \sum_{k=1}^{K} \left[g\left(\sum_{j=0}^{M} v_{kj} g\left(\sum_{i=0}^{d} w_{ji} x_{i} \right) \right) - t_{k} \right]^{2}$$

Apply chain rule:

$$\Delta w_{ji} = -\eta \frac{\partial E}{\partial w_{ji}} = -\eta \frac{\partial E}{\partial z_j} \frac{\partial z_j}{\partial a_j} \frac{\partial a_j}{\partial w_{ji}}$$
$$= -\eta \sum_{k=1}^{K} \underbrace{(y_k - t_k)g'(a_k)}_{\delta_k} v_{kj}g'(a_j)x_i$$
$$= -\eta \sum_{k=1}^{K} \delta_k v_{kj}g'(a_j)x_i$$

< 3 > < 3 >

Backprop. (Input-to-Hidden Layer) (cont.)

$$\Delta w_{ji} = -\eta \delta_j x_i$$

where

$$\delta_j = g'(a_j) \sum_{k=1}^{K} v_{kj} \delta_k$$

Observe: that we need to propagate the errors (δ 's) backwards to update the weights **v** and **w**

$$\begin{aligned} \Delta v_{kj} &= -\eta \delta_k z_j \\ \delta_k &= (y_k - t_k) g'(a_k) \\ \Delta w_{ji} &= -\eta \delta_j x_i \\ \delta_j &= g'(a_j) \sum_{k=1}^K v_{kj} \delta_k \end{aligned}$$

T.Stibor (TUM)

Course IN2207

SS2011 129 / 256

A B + A B +

Error Backpropagation

- Apply input **x** and forward propagate through the network using $a_j = \sum_{i=0}^{d} w_{ji}x_i$ and $z_j = g(a_j)$ to find the activations of all the hidden and output units.
- Compute the deltas δ_k for all the output units using $\delta_k = (y_k t_k)g'(a_k)$.
- Backpropagate the δ 's using $\delta_j = g'(a_j) \sum_{k=1}^{K} v_{kj} \delta_k$ to obtain δ_j for each hidden unit in the network.

Time and space complexity:

d input units, *M* hidden units and *K* output units results in M(d+1) weights in first layer and K(M+1) weights in second layer. Space and time complexity is $\mathcal{O}(M(K+d))$. If *e* training epochs are performed, then time complexity is $\mathcal{O}(e M(K+d))$.

< ロ > < 同 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Backprop. (Output-to-Hidden Layer) Vis.



★ E → E → Q へ C
 SS2011 131 / 256

Backprop. (Hidden-to-Input Layer) Vis.



★ E → E → Q Q Q
 SS2011 132 / 256

Property of Activation Functions

- In the Backpropagation algorithm the derivative of g(a) is required to evaluate the δ's.
- Activation functions

$$g_1(a) = rac{1}{1 + \exp(-eta a)}$$
 and $g_2(a) = anh(eta a)$

obey the property

$$egin{array}{rcl} g_1'(a) &=& eta \, g_1(a)(1-g_1(a)) \ g_2'(a) &=& eta(1-(g_2(a))^2) \end{array}$$

Online Backpropagation Algorithm

```
input : (\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N) \in \mathbb{R}^d \times \{C_1, C_2, \dots, C_K\}, \eta \in \mathbb{R}_+, \text{max.epoch} \in \mathbb{N}, \epsilon \in \mathbb{R}_+
output: w, v
begin
Randomly initialize w, v
epoch \leftarrow 0
repeat
for n \leftarrow 1 to N do
\mathbf{x} \leftarrow \text{select pattern } \mathbf{x}_n
v_{kj} \leftarrow v_{kj} - \eta \delta_k z_j
w_{ji} \leftarrow w_{ji} - \eta \delta_j x_i
epoch \leftarrow \text{epoch} + 1
until (epoch = max.epoch) or (\|\nabla E\| < \epsilon)
return w, v
```

< 回 > < 三 > < 三 >

Batch Backpropagation Algorithm

```
input : (\mathbf{x}_1, t_1), \ldots, (\mathbf{x}_N, t_N) \in \mathbb{R}^d \times \{\mathcal{C}_1, \mathcal{C}_2, \ldots, \mathcal{C}_K\}, \eta \in \mathbb{R}_+, \text{max.epoch} \in \mathbb{N}, \epsilon \in \mathbb{R}_+
output: w, v
begin
         Randomly initialize w, v
         epoch \leftarrow 0, \Delta w_{ji} \leftarrow 0, \Delta v_{ki} \leftarrow 0
         repeat
                   for n \leftarrow 1 to N do
                        \mathbf{x} \leftarrow \text{select pattern } \mathbf{x}_n
\Delta v_{kj} \leftarrow \Delta v_{kj} - \eta \delta_k z_j, \ \Delta w_{ji} \leftarrow \Delta w_{ji} - \eta \delta_j x_i
                   v_{kj} \leftarrow v_{kj} + \Delta v_{ki}
                   w_{ii} \leftarrow w_{ii} + \Delta w_{ii}
                   epoch \leftarrow epoch + 1
         until (epoch = max.epoch) or (\|\nabla E\| < \epsilon)
         return w. v
end
```

.

Multi-Layer Networks & Heaviside Step Func.



Possible decision boundaries which can be generated by networks having various numbers of layers and using Heaviside activation function.

/ -	
I Stubor (1 1 1 1 1 1
1.5000	1 0 1 1

Multi-Layer NN for XOR Separability Problem



<i>x</i> ₁	<i>x</i> ₂	$x_1 \text{ XOR } x_2$
-1	-1	-1
-1	+1	+1
+1	-1	+1
+1	+1	-1

$$g(a) = \left\{ egin{array}{cc} -1 & ext{if } a < 0 \ +1 & ext{if } a \geq 0 \end{array}
ight.$$

∃ >

Multi-Layer NN for XOR Sep. Problem (cont.)









T.Stibor (TUM)

Course IN2207

SS2011 138 / 256

Expressive Power of Multi-Layer Networks

With a two-layer network and a sufficient number of hidden units, *any* type of function can be represented when given proper nonlinearities and weights.

The famous mathematician Andrey Kolmogorov proved that any continuous function $y(\mathbf{x})$ defined on the unit hypercube $[0,1]^n$, $n \ge 2$ can be represented in the form

$$y(\mathbf{x}) = \sum_{j=1}^{2n+1} \Xi_j \left(\sum_{i=1}^d \Psi_{ij}(x_i) \right)$$

for properly chosen Ξ_j and Ψ_{ij} .

Bayes Decision Region vs. Neural Network



Points from blue and red class are generated by a mixture of Gaussians. Black curve shows optimal separation in a Bayes sense. Gray curve shows neural network separation of two independent backpropagation learning runs.

T.Stibor (TUM)

Neural Network (Density) Decision Region



SS2011 141 / 256

Overfitting/Underfitting & Generalization

Consider the problem of polynomial curve fitting where we shall fit the data using a polynomial function of the form:

$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \ldots + w_M x^M = \sum_{j=0}^M w_j x^j.$$

We measure the misfit of our predictive function $y(x, \mathbf{w})$ by means of error function which we like to minimize:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^{N} (y(x_i, \mathbf{w}) - t_i)^2$$

where t_i is the corresponding target value in the given training data set.

. .

Polynomial Curve Fitting



э 143 / 256 SS2011

-

A.

Polynomial Curve Fitting (cont.)

	M = 0	M = 1	<i>M</i> = 3	M = 9
w_0^*	0.19	0.82	0.31	0.35
w_1^{\star}		-1.27	7.99	232.37
w_2^{\star}			-25.43	-5321.83
w_3^{\star}			17.37	48568.31
w_4^{\star}				-231639.30
w_5^{\star}				640042.26
w_6^{\star}				-1061800.52
W_7^{\star}				1042400.18
w ₈ *				-557682.99
w_9^{\star}				125201.43

Table: Coefficients \mathbf{w}^* obtained from polynomials of various order. Observe the dramatically increase as the order of the polynomial increases (this table is taken from Bishop's book).

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >
Polynomial Curve Fitting (cont.)

Observe:

- if M is too small then the model underfits the data
- if M is too large then the model overfits the data

If M is too large then the model is more flexible and is becoming increasingly tuned to random noise on the target values. It is interesting to note that the overfitting problem become less severe as the size of the data set increases.



Polynomial Curve Fitting (cont.)

One technique that can be used to control the overfitting phenomenon is the *regularization*.

• Regularization involves adding a penalty term to the error function in order to discourage the coefficients from reaching large values.

The modified error function has the form:

$$\widehat{E}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^{N} (y(x_i, \mathbf{w}) - t_i)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}.$$

By means of the penalty term one reduces the value of the coefficients (shrinkage method).

Regularized Polynomial Curve Fitting M = 9



T.Stibor (TUM)

SS2011 147 / 256

Regularization in Neural Networks

- Number of input/output units is generally determined by the dimensionality of the data set.
- Number of hidden units *M* is free parameter that can be adjusted to obtain best predictive performance.
- Generalization error is not a simple function of *M* due to the presence of local minima in the error function.
- One straightforward way to deal with this problem is to increase stepwise the value of *M* and to choose the specific solution having the smallest test error.

Regularization in Neural Networks (cont.)

Equivalent to the regularized curve fitting approach, we can choose a relatively large value for M and control the complexity by the addition of a regularized term to the error function.

$$\widehat{E}(\mathbf{w}) = E(\mathbf{w}) + rac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

This form of regularization in neural networks is known as weight decay.

- Weight decay encourages weight values to decay towards zero, unless supported by the data.
- It can be considered as an example of a parameter shrinkage method because parameter values are shrunk towards zero.
- It can be also interpreted as the removal of non-useful connections during training.

イロト イポト イヨト イヨト

A too Overfitted Neural Network Model



Hidden units: 20, weight decay: 0

Hidden units: 20, weight decay: 0



イロト イヨト イヨト イヨト

A too Underfitted Neural Network Model



nite: 20 waight decay: 1

Hidden units: 20. weight decay: 2



<ロト < 団ト < 団ト < 団ト

Model Complexity is Properly Penalized



Hidden units: 20, weight decay: 0.3

Hidden units: 20. weight decay: 0.3



イロト イヨト イヨト イヨト

Regularization by Early Stopping

• Another alternative of regularization as a way of controlling the effective complexity of a network is the procedure of *early stopping*.



T.Stibor (TUM)

Course IN2207

SS2011 153 / 256

Example Early Stopping after 10 Epochs



Hidden units: 20, weight decay: 0, early stop after: 10

lidden units: 20, weight decay: 0, early stop after: 10



イロト イヨト イヨト イヨト

SS2011 154 / 256

Example Early Stopping after 50 Epochs



Hidden units: 20, weight decay: 0, early stop after: 50

ヘロト 人間ト 人造ト 人造ト

ht decay: 0, early stop after: 50

Example Early Stopping after 100 Epochs



Hidden units: 20, weight decay: 0, early stop after: 100



ヘロト 人間ト 人造ト 人造ト

idden units: 20, weight decay: 0, early stop after: 100

Cross-Validation & Prediction Error

Most commonly used method for estimating the prediction error (generalization error) on new data is *cross-validation*.



Split training data into K roughly equal-sized parts (e.g. K = 5). Learn neural network on K - 1 training parts and predict the not seen testing part. Perform this for k = 1, 2, ..., K and combine the K estimates of prediction error (test error).

T.Stibor (TUM)

Cross-Validation & Prediction Error (cont.)

Typical choices of K are 5 or 10. If training data set consists of N data points and K = N, then one obtains the *leave-one-out* cross-validation method.

```
# number of folds
folds <- 10
samps <- sample(rep(1:folds, length=n), n, replace=FALSE)</pre>
```

```
for (fold.iter in 1:folds) {
```

```
train <- dataset[samps!=fold.iter,] # fit the model
test <- dataset[samps==fold.iter,] # predict</pre>
```

perform here neural network training on set train
perform here neural network prediction on set test

}

・ 同 ト ・ ヨ ト ・ ヨ ト …

Bias-Variance Dilemma

- Idea: decompose generalization error into the two terms (bias,variance).
- Let us consider our curve fitting problem again. Let $h(\mathbf{x})$ be the true (but unknown) function with continuous valued output with noise.
- We want to estimate h(·) based on training data sets D each of size N.
- The natural way to measure the effectiveness of the estimator is to use the mean-square deviation from the desired optimal.
- \bullet Averaging over all training data sets ${\cal D}$ one gets the decomposition:

$$\mathbb{E}_{\mathcal{D}}\left[\{y(\mathbf{x}, \mathcal{D}) - h(\mathbf{x})\}^{2}\right]$$
$$= \underbrace{\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}, \mathcal{D})] - h(\mathbf{x})\}^{2}}_{\text{bias}^{2}} + \underbrace{\mathbb{E}_{\mathcal{D}}\left[\{y(\mathbf{x}, \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}, \mathcal{D})]\}^{2}\right]}_{\text{variance}}$$

Bias-Variance Dilemma (cont.)

The bias-variance dilemma is general phenomenon and also occurs as the under/overfitting problem in neural networks. In the context of neural networks:

- Bias is a measure of how much the network output, averaged over all possible data sets differs from the desired function.
- Variance is a measure of how much the network output varies between data sets.

In early stage of training, the bias is large because the network is far from the desired function. If the network is trained too long, then the network will also have learned the noise specific in the data set.

Bias-Variance Dilemma (cont.)



T.Stibor (TUM)

Course IN2207

SS2011 161 / 256

Momentum

- Gradient descent can be very slow if η is too small, and can oscillate widely if η is too large.
- Idea: use fraction of the previous weight change and actual gradient term to control non-radical revisions in the updates.

$$\Delta \mathbf{w}^{\mathsf{new}} = -\eta \frac{\partial E}{\partial \mathbf{w}} + \alpha \Delta \mathbf{w}, \qquad \mathbf{0} \le \alpha \le 1.$$

Momentum:

- can cancel side-to-side oscillations across the error valley,
- can cause a faster convergence when weight updates are all in the same direction because the learning rate is amplified.

Momentum Example Rosenbrock Function

Rosenbrock function $f(x, y) = (1 - x)^2 + 100(y - x^2)^2$ has global minimum f(x, y) = (0, 0) at (1, 1). Momentum param. $\alpha = 0.021$, learning rate $\eta = 0.001$.



Momentum Example Rosenbrock Func. (cont.) Setting $\alpha = 0$ (no momentum)



SS2011 164 / 256

∃ ▶ ∢

Momentum Example Rosenbrock Func. (cont.) Setting $\alpha = 0$ (no momentum) and a larger learning rate η



SS2011 165 / 256

Adaptive Parameter

- Adjust the parameters automatically as learning progress.
- Check if particular weight update decreased the error function.
 - If not we overleaped and η should be reduced.
 - \blacktriangleright If several steps in a row have decreased the error, we are too conservative and could try increasing η

$$\Delta\eta = \left\{ egin{array}{ll} +a & ext{if } \Delta E < 0 \ -b\eta & ext{if } \Delta E > 0 \ 0 & ext{otherwise} \end{array}
ight.$$

consistently (e.g. last K steps)

Hopfield Network Introductory Example

- Suppose we want to store N binary images in some memory.
- The memory should be content-addressable and insensitive to small errors.
- We present corrupted images to the memory (e.g. our brain) and recall the corresponding images.



presentation of corrupted images

Course IN220

recalled by the memory



- *w_{ij}* denotes weight connection from unit *j* to unit *i*
- no unit has connection with itself $w_{ii} = 0, \forall i$
- connections are symmetric
 w_{ij} = w_{ji}, ∀i, j

State of unit *i* can take values ± 1 and is denoted as S_i . State dynamics are governed by activity rule:

$$S_i = \operatorname{sgn}\left(\sum_j w_{ij}S_j\right), \text{ where } \operatorname{sgn}(a) = \begin{cases} +1 & \text{if } a \ge 0, \\ -1 & \text{if } a < 0 \end{cases}$$

Learning Rule in a Hopfield Network

Learning in Hopfield networks:

- Store a set of desired memories {x⁽ⁿ⁾} in the network, where each memory is a binary pattern with x_i ∈ {−1, +1}.
- The weights are set using the sum of outer products

$$w_{ij}=\frac{1}{N}\sum_{n}x_{i}^{(n)}x_{j}^{(n)},$$

where *N* denotes the number of units (*N* can also be some positive constant, e.g. number of patterns). Given a $m \times 1$ column vector **a** and $1 \times n$ row vector **b**. The outer product $\mathbf{a} \otimes \mathbf{b}$ (short $\mathbf{a} \mathbf{b}$) is defined as the $m \times n \max_{\substack{a \\ a_2 \\ a_2 \\ a_3 \\ a_3 \\ b_1 \\ a_3 \\ b_2 \\ a_3 \\ b_1 \\ a_3 \\ b_2 \\ a_3 \\ b_3 \\ b_3 \\ b_3 \\ b_3 \\ c_3 \\ b_3 \\ c_3 \\ b_3 \\ c_3 \\ c_3$

▲ロト ▲圖 ▶ ▲ 画 ▶ ▲ 画 ▶ ● の Q @

Learning in Hopfield Network (Example)

Suppose we want to store patterns $\mathbf{x}^{(1)} = [-1,+1,-1]$ and $\mathbf{x}^{(2)} = [+1,-1,+1].$

$$\begin{bmatrix} -1 \\ +1 \\ -1 \end{bmatrix} \otimes \begin{bmatrix} -1, +1, -1 \end{bmatrix} = \begin{bmatrix} +1 & -1 & +1 \\ -1 & +1 & -1 \\ +1 & -1 & +1 \end{bmatrix}$$
$$\begin{pmatrix} +1 \\ -1 \\ +1 & -1 & +1 \\ \end{bmatrix}$$
$$\begin{pmatrix} +1 \\ -1 \\ +1 & -1 \\ +1 & -1 \\ +1 & -1 & +1 \end{bmatrix}$$

SS2011 170 / 256

イロト イポト イヨト イヨト

Learning in Hopfield Network (Example) (cont.)

$$\mathbf{W} = \frac{1}{3} \begin{bmatrix} \mathbf{0} & -2 & +2 \\ -2 & \mathbf{0} & -2 \\ +2 & -2 & \mathbf{0} \end{bmatrix}$$

Recall: no unit has connection with itself.

The storage of patterns in the network can also be interpreted as constructing stable states. The condition for patterns to be stable is:

$$\operatorname{sgn}\left(\sum_{j} w_{ij} x_{i}\right) = x_{i}, \forall i.$$

Suppose we present pattern $\mathbf{x}^{(1)}$ to the network and want to restore the corresponding pattern.

Learning in Hopfield Network (Example) (cont.)

Let us assume that the network states are set as follows: $S_i = x_i$, $\forall i$. We can restore pattern $\mathbf{x}^{(1)} = [-1, +1, -1]$ as follows:

$$S_1 = \operatorname{sgn}\left(\sum_{j=1}^3 w_{1j}S_j\right) = -1 \qquad S_2 = \operatorname{sgn}\left(\sum_{j=1}^3 w_{2j}S_j\right) = +1$$
$$S_3 = \operatorname{sgn}\left(\sum_{j=1}^3 w_{3j}S_j\right) = -1$$

Can we also restore the original patterns by presenting "similar" patterns which are corrupted by noise?

Course IN2207

Updating States in a Hopfield Network

Synchronous updates:

• all units update their states $S_i = \text{sgn}\left(\sum_j w_{ij}S_j\right)$ simultaneously.

Asynchronous updates:

• one unit at a time updates its state. The sequence of selected units may be a fixed sequence or a random sequence.

Synchronously updating states can lead to oscillation (no convergence to a stable state).



A B K A B K

Aim of a Hopfield Network

Our aim is that by presenting a corrupted pattern, and by applying iteratively the state update rule the Hopfield network will settle down in a stable state which corresponds to the desired pattern.

Hopfield network is a method for

- pattern completion
- error correction.

The state of a Hopfield network can be expressed in terms of the energy function

$$E = -\frac{1}{2} \sum_{i,j} w_{ij} S_i S_j$$

Hopfield observed that if a state is a local minimum in the energy function, it is also a stable state for the network.

Basin of Attraction and Stable States



Within the space the stored patterns $\mathbf{x}^{(n)}$ are acting like attractors.

SS2011 175 / 256

Haykin's Digit Example

Suppose we stored the following digits in the Hopfield network:



< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Updated States of Corrupted Digit 6

Energy = -10.27



Energy = -12.2



Energy = -13.6



Energy = -14.87



Energy = -15.87



Energy = -18.07



Energy = -20.4







Energy = -22.2



Energy = -23.33



Energy = -25.73



Energy = -26.8



T.Stibor (TUM)

Energy = -29.67

updated unit 3



Energy = -30.13



Course IN2207



updated unit 6

Energy = -34.4



SS2011

177 / 256







Updated States of Corrupted Digit 6 (cont.)

Energy = -36.73



Energy = -38.4



Energy = -41.07



Energy = -42.4



Energy = -45.27



Energy = -47.6



Energy = -50.4



Energy = -52.67



Energy = -56.47



Energy = -58.4







Energy = -63.33



Energy = -64.47



Course IN2207

Energy = -68



Energy = -71.27



178 / 256

Updated States of Corrupted Digit 6 (cont.)

The resulting pattern (stable state with energy -90.47) matches the desired pattern.



Energy = -90.47







< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Recall a Spurious Pattern





Energy = -28.27





updated unit 12

Energy = -30.27

Energy = -31.93



Energy = -32.8



Energy = -33.4



Energy = -35.6



updated unit 111

Energy = -37.6





Energy = -42.6



Energy = -44.53



T.Stibor (TUM)

updated unit 15





Course IN2207

Energy = -50.53



updated unit 62









Energy = -40






Recall a Spurious Pattern (cont.)



Energy = -53.73

Energy = -56.53



Energy = -59.93



Energy = -61.6



Energy = -63.2



Energy = -63.73





Energy = -66.8

updated unit 112



updated unit 48





Energy = -71.93



T.Stibor (TUM)

Energy = -74.13

updated unit 70







Course IN2207

Energy = -80.27



Energy = -81.4



181 / 256



Energy = -69

Recall a Spurious Pattern (cont.)

The Hopfield network settled down in local minima with energy -84.93. This pattern however is *not* the desired pattern. It is a pattern which was not stored in the network.



A B K A B K

Incorrect Recall of Corrupted Pattern 2



Energy = -22.07

Energy = -22.07



Energy = -22.13



Energy = -22.33



Energy = -24.13



Energy = -24.53



Energy = -27.6



updated unit 100



Energy = -28.33

Energy = -29.87



Energy = -31.47



Energy = -32.13



T.Stibor (TUM)

Energy = -32.33

updated unit 86





Energy = -35.47

Course IN2207

Energy = -36.53



Energy = -38.67



183 / 256

Incorrect Recall of Corrupted Pattern 2 (cont.)

Energy = -39.2



Energy = -41.73







Energy = -48



Energy = -49.6



Energy = -51.6



Energy = -51.67



updated unit 37



Energy = -56.4



Energy = -58.27



Energy = -60.73



T.Stibor (TUM)

Energy = -61.87



Energy = -62.87



Course IN2207

Energy = -64.8



Energy = -68.93



184 / 256

Incorrect Recall of Corrupted Pattern 2 (cont.)



Energy = -69.87

updated unit 8

Energy = -70.13



Energy = -71.47



Energy = -72.93



Energy = -73.47



Energy = -77.07



Energy = -78.8



updated unit 67



updated unit 112

Energy = -82.67



Energy = -83.8







T.Stibor (TUM)





Energy = -86.4



Energy = -87.73



Energy = -88.53



Incorrect Recall of Corrupted Pattern 2 (cont.)

Although we presented the corrupted pattern 2, the Hopfield network settled down in the stable state that corresponds to pattern 6.



< 3 > < 3

MacKay's Example of an Overloaded Network

Six patterns are stored in the Hopfield network, however most of them are not stable states.



Spurious states represent stable states that are different from the stored desired patterns.

Spurious States and Capacity

- Reversed states $((-1) \cdot \mathbf{x}^{(n)})$ have same energy as the original patterns $\mathbf{x}^{(n)}$.
- Stable mixture states are not equal to any single pattern. They corresponds to a linear combination of an odd number of patterns.
- Spin glass states are local minima that are not correlated with any finite number of the original patterns.

Capacity:

What is the relation between the number d of units and the maximum number N_{max} of patterns one can store by allowing some small error. If $N_{\text{max}} = \frac{d}{4 \log d}$ then most of stored patterns can be recalled perfectly.

Support Vector Machine (SVM)



T.Stibor (TUM)

Course IN2207

SS2011 189 / 256

First papers on SVM

A Training Algorithm for Optimal Margin Classifiers, Bernhard E. Boser and Isabelle M. Guyon and Vladimir Vapnik, Proceedings of the fifth annual workshop on Computational learning theory (COLT), 1992 http://www.clopinet.com/isabelle/Papers/colt92.ps.Z

Support-Vector Networks, Corinna Cortes and Vladimir Vapnik, Machine Learning, 20(3):273-297, 1995 http://homepage.mac.com/corinnacortes/papers/SVM.ps

Extracting Support Data for a Given Task, Bernhard Schölkopf and Christopher J.C. Burges and Vladimir Vapnik, First International Conference on Knowledge Discovery & Data Mining, 1995 http://research.microsoft.com/~cburges/papers/kdd95.ps.gz

ヘロト 人間ト 人間ト 人間ト

Literature (Papers, Bookchapter, T.Report)

Statistical Learning and Kernel Method, Bernhard Schölkopf, Microsoft Research Technical Report, 2000, ftp://ftp.research.microsoft.com/pub/tr/tr-2000-23.pdf

An Introduction to Kernel-based Learning Algorithms, K.-R. Müller, S. Mika, G. Rätsch, K. Tsuda, and B. Schölkopf, IEEE Neural Networks, 12(2):181-201, 2001, http://ieeexplore.ieee.org/iel5/72/19749/00914517.pdf

Support Vector Machines, S. Mika, C. Schäfer, P. Laskov, D. Tax and K.-R. Müller, Bookchapter : Handbook of Computational Statistics (Concepts and Methods), http://www.xplore-stat.de/ebooks/scripts/csa/html/

Support Vector Machines for Classification and Regression, S.R. Gunn, Technical Report,

 $http://www.ecs.soton.ac.uk/{\sim}srg/publications/pdf/SVM.pdf$

ヘロト 人間ト 人造ト 人造ト

Literature (Paper, Books)

A Tutorial on Support Vector Machines for Pattern Recognition, Christopher J.C. Burges, Kluwer Academic Publishers, Boston http://research.microsoft.com/~cburges/papers/SVMTutorial.pdf

Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond, Bernhard Schölkopf and Alexander J. Smola, MIT Press, 2001

An Introduction to Support Vector Machines and Other Kernel-based Learning Methods, Nello Cristianini and John Shawe-Taylor, Cambridge University Press, 2000

Kernel Methods for Pattern Analysis, John Shawe-Taylor and Nello Cristianini, Cambridge University Press, 2004

The Nature of Statistical Learning Theory, Vladimir N. Vapnik, Springer-Verlag, sec. edition, 1999

イロト 不得 トイヨト イヨト

Literature (WWW and Source Code)

- http://www.kernel-machines.org/
- http://www.support-vector.net/
- http://www.learning-with-kernels.org/
- http://www.pascal-network.org/
- http://www.r-project.org/ (packages e1071,kernlab)

Scaling Freedom (Problem)



Multiplying **w** and *b* by the same constant κ gives the same hyperplane, represented in terms of different parameters

$$\kappa[(\mathbf{w}\cdot\mathbf{x})+b]=0\Leftrightarrow(\kappa\mathbf{w}\cdot\mathbf{x}+\kappa b)=0\Leftrightarrow(\mathbf{w}'\cdot\mathbf{x})+b'=0.$$

Example: $\mathbf{w} := (5, -4)^T$, b := 2; $5x - 4y + 2 = 0 \Leftrightarrow y = \frac{5}{4}x + \frac{1}{2}$, $\kappa := 3$; $15x - 12y + 6 = 0 \Leftrightarrow y = \frac{5}{4}x + \frac{1}{2}$

Canonical Hyperplane

A canonical hyperplane with respect to given m examples $(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_m, y_m) \in \mathcal{X} \times \{\pm 1\}$ is defined as a function

$$f(\mathbf{x}) = (\mathbf{w} \cdot \mathbf{x}) + b$$

where \mathbf{w} is normalized such that

$$\min_{i=1,\ldots,m}|f(\mathbf{x}_i)|=1$$

i.e. scaling **w** and *b* such that the points closest to the hyperplane satisfy $|(\mathbf{w}\cdot\mathbf{x}_i)+b|=1$

A B K A B K

Plus/Minus and Zero Hyperplane



SS2011 196 / 256

Optimal Separating Hyperplane



Let \mathbf{x}_3 be any point on the "minus" hyperplane and let \mathbf{x}_1 be the closest point to \mathbf{x}_3 .

$$\mathbf{w} \cdot \mathbf{x}_1 + b = +1$$
$$\mathbf{w} \cdot \mathbf{x}_3 + b = -1$$
$$\mathbf{x}_1 = \mathbf{x}_3 + \lambda \mathbf{w}$$
$$\mathbf{x}_1 \cdot \underbrace{(\mathbf{x}_3 + \lambda \mathbf{w})}_{\mathbf{x}_1} + b = 1$$

w

SS2011 197 / 256

Formulation as an Optimization Problem

Hyperplane with maximum margin (the smaller the norm of the weight vector \mathbf{w} , the larger the margin):

$$\begin{array}{ll} \text{minimize} & \displaystyle \frac{1}{2} \| \mathbf{w} \|^2 \\ \text{subject to} & \displaystyle y_i \left((\mathbf{w} \cdot \mathbf{x}_i) + b \right) \geq 1, \quad i = 1, \ldots, m \end{array}$$

Introduce Lagrange multipliers $\alpha_i \geq 0$ and a Lagrangian

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^m \alpha_i (y_i((\mathbf{x}_i \cdot \mathbf{w}) + b) - 1)$$

At the extrema, we have

$$rac{\partial}{\partial b}L(\mathbf{w},b,lpha)=0,\quad rac{\partial}{\partial \mathbf{w}}L(\mathbf{w},b,lpha)=0$$

Course IN2207

Optimization and Kuhn-Tucker Theorem

leads to solution

$$\sum_{i=1}^{m} \alpha_i y_i = 0, \qquad \mathbf{w} = \sum_{i=1}^{m} \alpha_i y_i \mathbf{x}_i.$$

The extreme point solution obtained has an important property that results from optimization known as the Kuhn-Tucker theorem. The theorem says:

Lagrange multiplier can be non-zero only if the corresponding inequality constraint is an equality at the solution.

This implies that *only* the training examples \mathbf{x}_i that lie on the plus and minus hyperplane have their corresponding α_i non-zero.

Relevant/Irrelevant Support Vector

More formally, the Karush-Kuhn-Tucker complementarity conditions say:

$$\alpha_i [y_i((\mathbf{x}_i \cdot \mathbf{w}) + b) - 1] = 0, \quad i = 1, \dots, m$$

the Support Vectors lie on the margin. That means for all training points

$$egin{array}{lll} [y_i((old x_i \cdot old w) + b)] > 1 & \Rightarrow lpha_i = 0 & o old x_i ext{ irrelevant} \ [y_i((old x_i \cdot old w) + b)] = 1 & (ext{on margin}) & o old x_i ext{ Support Vector} \end{array}$$

Relevant/Irrelevant Support Vector



Dual Form

The dual form has many advantages

- Formulate optimization problem without **w** (mapping **w** in high-dimensional spaces).
- Formulate optimization problem by means of α , y_i and dot product $\mathbf{x}_i \cdot \mathbf{x}_j$.
- Quadratic Programming Solver.

maximize
$$W(\alpha) = \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

subject to $\alpha_i \ge 0$, $i = 1, \dots, m$ and $\sum_{i=1}^{m} \alpha_i y_i = 0$

(3)

Hyperplane Decision Function

The solution is determined by the examples on the margin. Thus

$$f(\mathbf{x}) = \operatorname{sgn} \left((\mathbf{x} \cdot \mathbf{w}) + b \right)$$

= $\operatorname{sgn} \left(\sum_{i=1}^{m} y_i \alpha_i \left(\mathbf{x} \cdot \mathbf{x}_i \right) + b \right)$

where

$$\alpha_i [y_i((\mathbf{x}_i \cdot \mathbf{w}) + b) - 1] = 0, \quad i = 1, \dots, m$$

and

$$\mathbf{w} = \sum_{i=1}^{m} \alpha_i y_i \mathbf{x}_i$$

(4) (5) (4) (5)

Non-separable Case



Case where the constraints $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \ge 1$ cannot be satisfied, i.e. $\alpha_i \to \infty$

T.Stibor (TUM)

Course IN2207

SS2011 204 / 256

Relax Constraints (Soft Margin)

Modify the constraints to

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \ge 1 - \xi_i$$
, with $\xi_i \ge 0$

and add



i.e. distance of error points to their correct place in the objective function

minimize
$$\frac{1}{2} \|\mathbf{w}\|^2 + C \cdot \sum_{i=1}^m \xi_i$$

Same dual, with additional constraints $0 \le \alpha_i \le C$

Course IN2207

Non-separable Case (Part 2)



This data set is not properly separable with lines (also when using many slack variables)

Separate in Higher-Dim. Space

Map data in higher-dimensional space and separate it there with a hyperplane



Feature Space

Apply the mapping

$$egin{array}{rcl} \Phi : \mathbb{R}^{N} &
ightarrow & \mathcal{F} \ \mathbf{x} & \mapsto & \Phi(\mathbf{x}) \end{array}$$

to the data $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m \in \mathcal{X}$ and construct separating hyperplane in \mathcal{F} instead of \mathcal{X} . The samples are preprocessed as $(\Phi(\mathbf{x}_1), y_1), \dots, (\Phi(\mathbf{x}_m), y_m) \in \mathcal{F} \times \{\pm 1\}.$

Obtained decision function:

$$\begin{aligned} F(\mathbf{x}) &= \operatorname{sgn}\left(\sum_{i=1}^{m} y_i \alpha_i \left(\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}_i)\right) + b\right) \\ &= \operatorname{sgn}\left(\sum_{i=1}^{m} y_i \alpha_i k(\mathbf{x}, \mathbf{x}_i) + b\right) \end{aligned}$$

How about patters $\mathbf{x} \in \mathbb{R}^N$ and product features of order d? Dim (\mathcal{F}) grows like N^d . Example $N = 16 \times 16$, and $d = 5 \longrightarrow \text{dimension} 10^{10}$

T.Stibor (TUM)

SS2011 208 / 256

Kernels

A *kernel* is a function k, such that for all $\mathbf{x}, \mathbf{y} \in \mathcal{X}$

$$k(\mathbf{x}, \mathbf{y}) = (\Phi(\mathbf{x}) \cdot \Phi(\mathbf{y})),$$

where Φ is a mapping from \mathcal{X} to an dot product feature space \mathcal{F} .

The $m \times m$ matrix K with elements $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ is called kernel matrix or Gram matrix. The kernel matrix is symmetric and positive semi-definite, i.e. for all $a_i \in \mathbb{R}$, i = 1, ..., m, we have $\sum_{i,j=1}^{m} a_i a_j K_{ij} \ge 0$

Positive semi-definite kernels are exactly those giving rise to a positive semi-definite kernel matrix K for all m and all sets $\{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_m\} \subseteq \mathcal{X}$.

The Kernel Trick Example

Example : compute 2nd order products of two "pixels", i.e.

$$\mathbf{x} = (x_1, x_2)$$
 and $\Phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$

$$\begin{aligned} (\Phi(\mathbf{x}) \cdot \Phi(\mathbf{z})) &= (x_1^2, \sqrt{2}x_1x_2, x_2^2)(z_1^2, \sqrt{2}z_1z_2, z_2^2)^T \\ &= ((x_1, x_2)(z_1, z_2)^T)^2 \\ &= (\mathbf{x} \cdot \mathbf{z}^T)^2 \\ &= : k(\mathbf{x}, \mathbf{z}) \end{aligned}$$

Kernel without knowing Φ

Recall: mapping $\Phi : \mathbb{R}^N \to \mathcal{F}$. SVM depends on the data through dot products in \mathcal{F} , i.e. functions of the form

 $\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$

With k such that k(x_i, x_j) = Φ(x_i) · Φ(x_j), it is not necessary to even know what Φ(x) is.

Example: $k(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u}-\mathbf{v}\|^2}{\gamma}\right)$, in this example \mathcal{F} is infinite dimensional.

ヘロト 人間ト 人間ト 人目ト

Feature Space (Optimization Problem)

Quadratic optimization problem (soft margin) with kernel:

maximize
$$W(\alpha) = \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$$

subject to $0 \le \alpha_i \le C$, $i = 1, \dots, m$ and $\sum_{i=1}^{m} \alpha_i y_i = 0$

(Standard) Kernels

Linear
$$k_0(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})$$

Polynomial $k_1(\mathbf{u}, \mathbf{v}) = ((\mathbf{u} \cdot \mathbf{v}) + \Theta)^d$
Gaussian $k_2(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|^2}{\gamma}\right)$
Sigmoidal $k_3(\mathbf{u}, \mathbf{v}) = \tanh(\kappa(\mathbf{u} \cdot \mathbf{v}) + \Theta)$

Course IN2207

・ロト ・ 日 ト ・ ヨ ト ・ ヨ ト

SVM Results for Gaussian Kernel





$$\gamma=$$
 0.5, $C=$ 50

$$\gamma = 0.5, \ C = 1$$

T.Stibor (TUM)

Course IN2207

SS2011 214 / 256

SVM Results for Gaussian Kernel (cont.)





$$\gamma = 0.02, \ C = 50$$

$$\gamma = 10, \ C = 50$$

T C	
I Stibor I	
1.5000	10111

Course IN2207

SS2011 215 / 256

Link to Statistical Learning Theory

Pattern Classification:

Learn $f: \mathcal{X} \rightarrow \{\pm 1\}$ from input-output training data

We assume that data is generated from some unknown (but fixed) probability distribution $P(\mathbf{x}, y)$


(Empirical) Risk

Learn f from training set $(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_m, y_m) \in \mathcal{X} \times \{\pm 1\}$, such that the expected missclassification error on a test set, also drawn from $P(\mathbf{x}, y)$,

$$R[f] = \int \frac{1}{2} |f(\mathbf{x}) - y| dP(\mathbf{x}, y)$$
 Expected Risk

is minimal.

Problem: we cannot minimize the expected risk, because P is unknown. Minimize instead the average risk over training set, i.e.

$$R_{\text{emp}}[f] = \frac{1}{m} \sum_{i=1}^{m} \frac{1}{2} |f(\mathbf{x}_i) - y_i|$$
 Empirical Risk

Problems with Empirical Risk

Minimizing the empirical risk (training error), does not imply a small test error. To see this, note that for each function f and any test set $(\overline{\mathbf{x}}_1, \overline{y}_1), \ldots, (\overline{\mathbf{x}}_m, \overline{y}_m) \in \mathcal{X} \times \{\pm 1\}$, satisfying $\{\overline{\mathbf{x}}_1, \ldots, \overline{\mathbf{x}}_{\overline{m}}\} \cap \{\mathbf{x}_1, \ldots, \mathbf{x}_m\} = \{\}$, there exists another function f^* such that $f^*(\mathbf{x}_i) = f(\mathbf{x}_i)$ for all $i = 1, \ldots, m$, but $f^*(\overline{\mathbf{x}}_i) \neq f(\overline{\mathbf{x}}_i)$ (on all test samples) for all $i = 1, \ldots, \overline{m}$



A Bound for Pattern Classification

For any $f \in F$ and m > h, with a probability of at least $1 - \eta$,



T.Stibor (TUM)

Course IN2207

SS2011 219 / 256

Measure Complexity of Function Set

How to measure the complexity of a given function set ?

The Vapnik-Chervonenkis (short VC) dimension is defined as the maximum number of points that can be labeled in all possible ways.

- In \mathbb{R}^2 we can shatter three non-collinear points.
- But we can never shatter four points in \mathbb{R}^2 .
- Hence the VC dimension is h = 3.



VC Dimension

- Separating hyperplanes in \mathbb{R}^N have VC dimension N+1.
- Hence: separating hyperplanes in high-dimensional feature spaces have extremely large VC dimension, and may not generalize well.
- But, "margin" hyperplanes can still have a small VC dimension.

VC Dimension of Margin Hyperplanes

Consider hyperplanes $(\mathbf{w} \cdot \mathbf{x}) = 0$ where w is normalized such they are in a canonical form w.r.t. a set of points $X^* = {\mathbf{x}_1, \ldots, \mathbf{x}_r}$, i.e.,

$$\min_{i=1,\ldots,r} |(\mathbf{w} \cdot \mathbf{x}_i)| = 1$$

The set of decision functions $f_{\mathbf{w}}(\mathbf{x}) = \operatorname{sgn}(\mathbf{x} \cdot \mathbf{w})$ defined on X^* and satisfying the constraint $\|\mathbf{w}\| \leq \Lambda$ has VC dimension satisfying

$$h \leq \min(R^2 \Lambda^2 + 1, N + 1)$$

Here, R is the smallest sphere around the origin containing X^* .

Smallest Sphere and VC Dimension



Hyperplanes with a large margin $(\frac{2}{\Lambda_1})$ induce only a small number of possibilities to separate the data, i.e. the VC dimension is small (left figure). In contrast, smaller margins $(\frac{2}{\Lambda_2})$ induce more separating hyperplanes, i.e. the VC dimension is large (right figure).

SVM (RBF Kernel) and Bayes Decision





< ∃ >

SS2011 224 / 256

Implementation (Matrix notation)

Recall:

maximize
$$W(\alpha) = \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

subject to $\alpha_i \ge 0$, $i = 1, \dots, m$ and $\sum_{i=1}^{m} \alpha_i y_i = 0$

This can be expressed in a matrix notation as

$$\min_{\alpha} \frac{1}{2} \alpha^{T} H \alpha + c^{T} \alpha$$

T.Stibor (TUM)

SS2011 225 / 256

(3)

Implementation (Matrix notation) cont.

where

$$H = ZZ^T, \quad c^T = (-1, \ldots, -1)$$

with constraints

$$\alpha^T Y = 0, \alpha_i \ge 0, i = 1, \dots, m$$

where

$$Z = \begin{bmatrix} y_1 \mathbf{x}_1 \\ y_2 \mathbf{x}_2 \\ \vdots \\ y_m \mathbf{x}_m \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}.$$

(3)

Implementation (Maple)

CalcAlphaVector := proc(LPM::Matrix, Kernel, C)

```
rows := RowDimension(LPM);
H := Matrix(1..rows,1..rows);
a := Vector(2);
b := Vector(2);
beq := Vector([0]);
Y := Vector(rows);
for k from 1 to rows do
 Y[k] := LPM[k,3];
end do;
```

A B K A B K

Implementation (Maple) cont.

```
for i from 1 to rows do
  for j from 1 to rows do
    a[1] := LPM[i,1];
    a[2] := LPM[i,2];
    b[1] := LPM[j,1];
    b[2] := LPM[j,2];
   H[i,j] := Y[i] * Y[j] * Kernel(a,b);
  end do;
end do;
```

▲ E → E → Q < C SS2011 228 / 256

A B + A B +

Implementation (Maple) cont.

```
c := Vector(1..rows, fill = -1);
```

```
bl := Vector(1..rows, fill = 0);
```

```
bu := Vector(1..rows, fill = C);
```

```
Aeq := convert(Transpose(Y),Matrix);
```

return QPSolution[2];

end proc:

A B + A B +

```
Using SVM in R (e1071)
```

Example, how to use SVM classification in R:

```
library(e1071);
```

```
# load famous Iris Fisher data set
data(iris);
attach(iris);
```

A B + A B +

Using SVM in R (e1071) cont.

```
# Using the first fold:
train <- iris[samps!=1,] # fit the model
test <- iris[samps==1,] # predict</pre>
```

```
# features 1 -> 4 from training set
x_train <- train[,1:4];
# class labels
y_train <- train[,5];</pre>
```

Using SVM in R (e1071) cont.

```
# features 1 -> 4 from test set (no class labels)
x_test <- test[,1:4];</pre>
```

```
# predict class labels for x_test
pred <- predict(model,x_test,decision.values = TRUE);</pre>
```

```
# get true class labels
y_true <- test[,5];</pre>
```

```
# check accuracy:
table(pred, y_true);
# determine test error
sum(pred != y_true)/length(y_true);
# do that for all folds, i.e. train <- iris[samps!=2,]
# test <- iris[samps==2,] , .....</pre>
```

◆□▶ ◆□▶ ◆三▶ ◆三▶ ● ○ ○ ○

SVM Multi-Class Classification

- A SVM is a binary classifier, that is, the class labels can only take two values: ±1.
- Many real-world problems, however, have more than two classes (e.g. optical character recognition).

One Versus the Rest: To get *M*-class classifiers, construct set of binary classifiers f^1, f^2, \ldots, f^M , each trained to separate one class from rest.

Combine them to get a multi-class classification according to the maximal output *before* applying the sgn function.

$$\underset{j=1...M}{\operatorname{argmax}} g^{j}(x), \text{ where } g^{j}(x) = \sum_{i=1}^{m} y_{i} \alpha_{i}^{j} k(x, x_{i}) + b^{j}.$$

$$f^{j}(x) = \operatorname{sgn}(g^{j}(x))$$

and

SVM Multi-Class Classification (cont.)

- Recall: $g^{j}(x)$ returns a signed real-valued value which can be interpreted as the distance from the separation (hyper)plane to the point x.
- Value can also be interpreted as a confidence value. The larger the value the more *confident* one is that the point *x* belong to the positive class.
- Hence, assign point x to the class whose confidence value is largest for this point.

SVM Pairwise Classification

- Train a classifier for each possible pair of classes.
- For *M* classes, this results in $\binom{M}{2} = \frac{(M-1)M}{2}$ binary classifiers.
- Classify an unknown point x by applying each of the $\binom{M}{2}$ binary classifiers and count how many times point x was assigned to that class label.
- Class label with highest count is then the considered label for the unknown point *x*.

One-Class SVM for Novelty Detection

Idea: enclose data with a hypersphere and classify new data as *normal* if it falls within the hypersphere and otherwise as anomalous data.



Course IN2207

Minimum Enclosing Hypersphere

Given normal data $\mathcal{X} = {\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m} \in \mathbb{R}^d$ and let r be the radius of the hypersphere and $\mathbf{c} \in \mathbb{R}^d$ the center. To find the minimum enclosing hypersphere we have to solve the following optimization problem:

minimize
$$r^2$$

subject to $\|\Phi(\mathbf{x}_i) - \mathbf{c}\|^2 \le r^2, \quad i = 1, \dots, m.$

Lagrangian multiplier $\alpha_i \geq 0$ for each constraint

$$L(\mathbf{c}, r, \alpha) = r^2 + \sum_{i=1}^m \alpha_i \left\{ \|\Phi(\mathbf{x}_i) - \mathbf{c}\|^2 - r^2 \right\}$$

T.Stibor (TUM)

Setting the derivatives with respect to \mathbf{c} and r to zero

$$\frac{\partial L(\mathbf{c}, r, \alpha)}{\partial \mathbf{c}} = 2\sum_{i=1}^{m} \alpha_i (\Phi(\mathbf{x}_i) - \mathbf{c}) = \mathbf{0}$$
$$\frac{\partial L(\mathbf{c}, r, \alpha)}{\partial r} = 2r \left(1 - \sum_{i=1}^{m} \alpha_i\right) = 0$$

one obtains the following equations

$$\sum_{i=1}^{m} \alpha_i = 1 \text{ and } \mathbf{c} = \sum_{i=1}^{m} \alpha_i \Phi(\mathbf{x}_i).$$
(2)

T.Stibor (TUM)

SS2011 238 / 256

Inserting relation (2) into

$$L(\mathbf{c}, r, \alpha) = r^2 + \sum_{i=1}^m \alpha_i \left\{ \|\Phi(\mathbf{x}_i) - \mathbf{c}\|^2 - r^2 \right\}$$
$$= \sum_{i=1}^m \alpha_i \|\Phi(\mathbf{x}_i) - \mathbf{c}\|^2$$
$$= \sum_{i=1}^m \alpha_i k(\mathbf{x}_i, \mathbf{x}_i) - \sum_{i,j=1}^m \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j)$$

gives the dual form.

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

To find α in dual form, solve optimization problem:

maximize
$$W(\alpha) = \sum_{i=1}^{m} \alpha_i k(\mathbf{x}_i, \mathbf{x}_i) - \sum_{i,j=1}^{m} \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j)$$

subject to $\sum_{i=1}^{m} \alpha_i = 1$, and $\alpha_i \ge 0$, $i = 1, \dots, m$.

Recall: Lagrange multiplier can be non-zero only if the corresponding inequality constraint is an equality at the solution.

The KKT complementarity conditions are satisfied by the optimal solutions α , (**c**, *r*)

$$\alpha_i \left\{ \| \Phi(\mathbf{x}_i) - \mathbf{c} \|^2 - r^2 \right\}, \quad i = 1, \dots, m.$$

This implies that only training examples \mathbf{x}_i that lie on the surface of the optimal hypersphere have their corresponding $\alpha_i > 0$.



T.Stibor (TUM)

Course IN2207

SS2011 241 / 256

Decision Function

$$f(\mathbf{x}) = \operatorname{sgn}(r^2 - \|\Phi(\mathbf{x}) - \mathbf{c}\|^2)$$

= $\operatorname{sgn}\left(r^2 - \left\{(\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x})) - 2\sum_{i=1}^m \alpha_i(\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}_i)) + \sum_{i,j=1}^m \alpha_i \alpha_j(\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j))\right\}\right)$
= $\operatorname{sgn}\left(r^2 - \left\{k(\mathbf{x}, \mathbf{x}) - 2\sum_{i=1}^m \alpha_i k(\mathbf{x}, \mathbf{x}_i) + \sum_{i,j=1}^m \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j)\right\}\right)$

T.Stibor (TUM)

▲ ■ → ■ → Q へ ⊂ SS2011 242 / 256

▲□ > ▲圖 > ▲ 圖 > ▲ 圖 >

Soft Enclosing Hypersphere

If we have some noise in our training set the "hard" enclosing hypersphere approach may force a larger radius than should really be needed. In other words, the solution would not be *robust*.

Aim: Find minimum enclosing hypersphere that contains (allmost) all training examples, but not some small portion of extreme training examples.



Soft Enclosing Hypersphere (cont.)

Introduce slack variables $\boldsymbol{\xi}, \xi_i \geq 0, i = 1, \dots, m$

$$\begin{array}{ll} \text{minimize} & r^2 + C \sum_{i=1}^m \xi_i \\ \text{subject to} & \|\Phi(\mathbf{x}_i) - \mathbf{c}\|^2 \leq r^2 + \xi_i, \quad \xi_i \geq 0, \ i = 1, \dots, m. \end{array}$$

Lagrangian multiplier $\alpha_i, \beta_i \geq 0$ for each constraint

$$L(\mathbf{c}, r, \boldsymbol{\alpha}, \boldsymbol{\beta}) = r^2 + C \sum_{i=1}^m \xi_i + \sum_{i=1}^m \alpha_i \left\{ \|\Phi(\mathbf{x}_i) - \mathbf{c}\|^2 - r^2 - \xi_i \right\} - \sum_{i=1}^m \beta_i \xi_i$$

T.Stibor (TUM)

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Soft Enclosing Hypersphere (cont.)

Setting partial derivatives to $\mathbf{0}$ gives

$$\sum_{i=1}^{m} \alpha_i = 1, \quad \mathbf{c} = \sum_{i=1}^{m} \alpha_i \Phi(\mathbf{x}_i)$$

This leads to the dual form

$$\begin{array}{ll} \text{minimize} & \sum_{i,j=1}^{m} \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^{m} \alpha_i k(\mathbf{x}_i, \mathbf{x}_i) \\ \text{subject to} & 0 \le \alpha_i \le C, \quad \sum_{i=1}^{m} \alpha_i = 1 \end{array}$$

Course IN2207

< ∃ > <

Hyperplane One-Class SVM

Idea: Separate in high-dimensional feature space \mathcal{F} , the points from the origin (circled point) with a maximum distance, and allow $\nu \cdot m$ many "outliers" which lie between the origin and the hyperplane, i.e. the -1 side.



T.Stibor (TUM)

Hyperplane One-Class SVM (cont.)

Normal vector of the hyperplane is determined by solving the primal quadratic optimization problem

minimize
$$\frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{\nu m} \sum_i \xi_i - \rho \tag{3}$$

subject to
$$(\mathbf{w} \cdot \Phi(\mathbf{x}_i)) \ge \rho - \xi_i, \xi_i > 0, i = 1, \dots, m.$$
 (4)

Lagrangian multiplier $\alpha_i, \beta_i \geq 0$ for each constraint

$$L(\mathbf{w}, \boldsymbol{\xi}, \rho, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{\nu m} \sum_i \xi_i - \rho$$
$$- \sum_{i=1}^m \alpha_i ((\mathbf{w} \cdot \Phi(\mathbf{x}_i)) - \rho + \xi_i) - \sum_{i=1}^m \beta_i \xi_i$$

Reformulating (3) and (4) to a dual optimization problem in terms of a kernel function $k(\cdot, \cdot)$, one obtains

Hyperplane One-Class SVM (cont.)

maximize
$$\frac{1}{2} \sum_{i,j=1}^{m} \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j)$$
 (5)

subject to
$$0 \le \alpha_i \le \frac{1}{\nu m}, i = 1, \dots, m$$
 and $\sum_{i=1}^m \alpha_i = 1.$ (6)

Differentiating the primal with respect to **w**, one gets $\mathbf{w} = \sum_{i=1}^{m} \alpha_i \Phi(\mathbf{x}_i)$.

Recall KKT theorem: For $\alpha_i > 0$ the corresponding pattern \mathbf{x}_i satisfies

$$\rho = (\mathbf{w} \cdot \Phi(\mathbf{x}_i)) = \sum_{j=1}^m \alpha_j k(x_j, x_i)$$

Hyperplane One-Class SVM (cont.)

The decision function (left/right side of the hyperplane):

$$f(\mathbf{x}) = \operatorname{sgn}((\mathbf{w} \cdot \Phi(\mathbf{x}_i)) - \rho)$$

= $\operatorname{sgn}\left(\sum_{i=1}^m \alpha_i k(\mathbf{x}_i, \mathbf{x}) - \rho\right)$

ν-Property:

- ν is an upper bound on the fraction of outliers.
- ν is a lower bound on the fraction of Support Vectors.

Hyperplane One-Class SVM Example



$$\nu = 0.05$$

 $\nu = 0.5$

イロト イヨト イヨト イヨト

	/ · · · · · · · ·
1 Stubor	
1.5000	

Course IN220

★ E → E → Q ←
 SS2011 250 / 256

Support Vector Regression

Basic idea: map the data **x** into a high-dimensional feature space \mathcal{F} via a nonlinear mapping Φ , and do *linear* regression in this space.

$$f(\mathbf{x}) = (\mathbf{w} \cdot \Phi(\mathbf{x})) + b$$
 with $\Phi : \mathbb{R}^d \to \mathcal{F}, \mathbf{w} \in \mathcal{F}$.

Linear regression in a high dimensional feature space corresponds to *nonlinear* regression in the low dimensional space \mathbb{R}^d .

Vapnik's ϵ -insensitive loss function:

$$|y - f(\mathbf{x})|_{\epsilon} := \max\{0, |y - f(\mathbf{x})| - \epsilon\}$$

Find function $f(\mathbf{x})$ that has at most ϵ deviation from all the targets y_i

Support Vector Regression (cont.)

Estimate linear regression $f(\mathbf{x}) = (\mathbf{w} \cdot \Phi(\mathbf{x})) + b$ leads to the problem of minimizing the term

$$\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n |y_i - f(\mathbf{x}_i)|_{\epsilon}$$

In the soft margin case one needs two types of slack variables $(\boldsymbol{\xi}, \boldsymbol{\xi}^*)$ for the two cases $f(\mathbf{x}_i) - y_i > \epsilon$ and $y_i - f(\mathbf{x}_i) > \epsilon$.



Figure is taken from Schölkopf's and Smola's book (Learning with Kernels)

T.Stibor (TUM)

Course IN2207

SS2011 252 / 256
Support Vector Regression (cont.)

Optimization problem is given by:

$$\begin{array}{ll} \text{minimize} & \frac{1}{2} \|\mathbf{w}\|^2 + C \cdot \sum_{i=1}^n (\xi_i + \xi_i^*) \\ \text{subject to} & f(\mathbf{x}_i) - y_i &\leq \epsilon + \xi_i \\ y_i - f(\mathbf{x}_i) &\leq \epsilon + \xi_i^* \\ \xi_i, \xi_i^* \geq 0 & \text{for all } i = 1, \dots, n \end{array}$$

A B K A B K

Support Vector Regression (cont.)

Introducing Lagrange multipliers α, α^* (dual form):

maximize
$$-\epsilon \sum_{i=1}^{n} (\alpha_{i}^{*} + \alpha_{i}) + \sum_{i=1}^{n} (\alpha_{i}^{*} - \alpha_{i})y_{i}$$
$$-\frac{1}{2} \sum_{i,j}^{n} (\alpha_{i}^{*} - \alpha_{i})(\alpha_{j}^{*} - \alpha_{j})k(\mathbf{x}_{i}, \mathbf{x}_{j})$$
subject to
$$0 \le \alpha_{i}, \alpha_{i}^{*} \le C \text{ for all } i = 1, \dots, n \text{ and}$$
$$\sum_{i=1}^{n} (\alpha_{i}^{*} - \alpha_{i}) = 0$$

Regression estimate takes the form

$$f(\mathbf{x}) = \sum_{i=1}^{n} (\alpha_i^* - \alpha_i) k(\mathbf{x}_i, \mathbf{x}) + b$$

(3)

Support Vector Regression (cont.)

Offset *b* can be computed by exploiting Karush-Kuhn-Tucker conditions: $f(\mathbf{x}_i) - y_i \le \epsilon + \xi_i$ becomes an equality with $\xi_i = 0$ if $0 < \alpha_i < C$ and $y_i - f(\mathbf{x}_i) \le \epsilon + \xi_i^*$ becomes an equality with $\xi_i^* = 0$ if $0 < \alpha_i^* < C$ that is:

$$\begin{aligned} \alpha_i(\epsilon + \xi_i - y_i + (\mathbf{w} \cdot \Phi(\mathbf{x}_i)) + b) &= 0 \\ \alpha_i^*(\epsilon + \xi_i^* + y_i - (\mathbf{w} \cdot \Phi(\mathbf{x}_i)) - b) &= 0 \end{aligned}$$

and leads to solution

$$b = y_i - (\mathbf{w} \cdot \Phi(\mathbf{x}_i)) - \epsilon \quad \text{for } \alpha_i \in (0, C)$$

$$b = y_i - (\mathbf{w} \cdot \Phi(\mathbf{x}_i)) + \epsilon \quad \text{for } \alpha_i^* \in (0, C)$$

T.Stibor (TUM)

Course IN2207

SS2011 255 / 256

Support Vector Regression Example



```
library(kernlab);
x <- seq(-20,20,0.1);
y <- sin(x)/x + rnorm(401,sd=0.03);
# train SVM
reg.svm <- ksvm(x,y,epsilon=0.01,kpar=list(sigma=16),cross=3);
plot(x,y,type="1",lwd=3);
lines(x,predict(reg.svm,x),col="red",lwd=3);
```

< ロト < 同ト < ヨト < ヨト