

Using Support Vector Machines


Machine Learning Tutorial for Supervised Classification

Dr. – Ing. Morris Riedel
 Adjunct Associated Professor
 School of Engineering and Natural Sciences, University of Iceland
 Research Group Leader, Juelich Supercomputing Centre, Germany


LECTURE 3

Support Vector Machines

July 7th, 2016
 Barcelona



UNIVERSITY OF ICELAND
 SCHOOL OF ENGINEERING AND NATURAL SCIENCES
 FACULTY OF INDUSTRIAL ENGINEERING,
 MECHANICAL ENGINEERING AND COMPUTER SCIENCE



JÜLICH
 FORSCHUNGSZENTRUM

Review of Lecture 2

- Supervised Classification
 - Finding a **target distribution** for realistic learning situations
 - Assume **unknown probability distribution** over the input space
 - Hypothesis search** with M models and we pick one
- Statistical Learning Theory

Unknown Target Distribution: $P(y|x)$
 target function $f: X \rightarrow Y$ plus noise
 (ideal function)

$x = (x_1, \dots, x_d)$ ← Probability Distribution P on X

Hypothesis Set
 $\mathcal{H} = \{h\}; g \in \mathcal{H}$
(set of candidate formulas)

Final Hypothesis
 $g \approx f$

$\Pr [| E_{in}(g) - E_{out}(g) | > \epsilon] \leq 2Me^{-2\epsilon^2 N}$

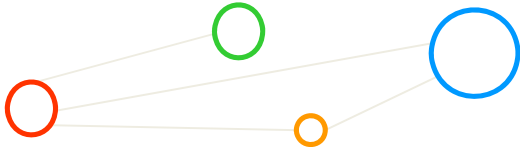
Learning Algorithm ('train a system') \mathcal{A} (set of known algorithms) → 'learn: get error smaller' $E_{in}(g) \approx 0$ → 'generalize well for unseen data' $E_{out}(g) \approx E_{in}(g)$

$\Pr [| E_{in}(g) - E_{out}(g) | > \epsilon] \leq 4m_{\mathcal{H}}(2N)e^{-1/8\epsilon^2 N}$

Shift the view to the data and we exchange M with growth function that is indeed depending on N

Lecture 3 – Support Vector Machines 2 / 58

Outline



Lecture 3 – Support Vector Machines 3 / 58

Outline of the Course

1. Machine Learning Fundamentals
2. Supervised Classification
3. Support Vector Machines
4. Applications and Serial Computing Limits
5. Kernel Methods
6. Applications and Parallel Computing Benefits

Lecture 3 – Support Vector Machines

4 / 58

Outline

- Maximal Margin Classifier
 - Term Support Vector Machines Refined
 - Margin as Geometric Interpretation
 - Optimization Problem & Implementation
 - Solving and Limitations of Classifier
 - Apply Classifier to Flower Problem
- Support Vector Classifier
 - From Hard-margin to Soft-margin
 - Understand the slack variables
 - Role of Regularization Parameter C
 - Solving and Limitations of Classifier
 - Apply Classifier to Flower Problem

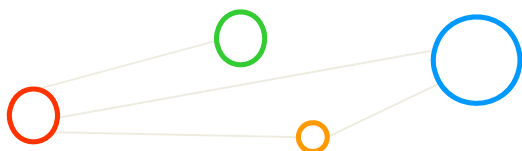
▪ Promises from previous lectures reviewed...
▪ Lecture 2: Lecture 3 provides details on how support vector machines benefit from regularization methods



Lecture 3 – Support Vector Machines

5 / 58

Maximum Margin Classifier

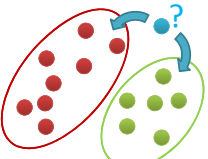
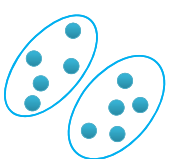
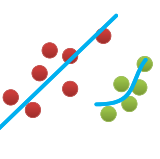


Lecture 3 – Support Vector Machines

6 / 58

Methods Overview – Focus in this Lecture

- Statistical data mining methods can be roughly categorized in classification, clustering, or regression augmented with various techniques for data exploration, selection, or reduction

Classification	Clustering	Regression
		
<ul style="list-style-type: none"> Groups of data exist New data classified to existing groups 	<ul style="list-style-type: none"> No groups of data exist Create groups from data close to each other 	<ul style="list-style-type: none"> Identify a line with a certain slope describing the data

Lecture 3 – Support Vector Machines 7 / 58

Term Support Vector Machines Refined

- Support Vector Machines (SVMs) are a classification technique developed ~1990
- SVMs perform well in many settings & are considered as one of the best 'out of the box classifiers'

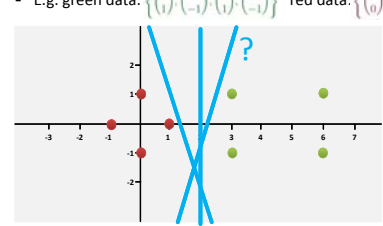
[1] An Introduction to Statistical Learning

- Term detailed refinement into 'three separate techniques'
 - Practice: applications mostly use the SVMs with kernel methods
- 'Maximal margin classifier'
 - A simple and intuitive classifier with a 'best' linear class boundary
 - Requires that data is 'linearly separable'
- 'Support Vector Classifier'
 - Extension to the maximal margin classifier for non-linearly separable data
 - Applied to a broader range of cases, idea of 'allowing some error'
- 'Support Vector Machines' → Using Non-Linear Kernel Methods
 - Extension of the support vector classifier
 - Enables non-linear class boundaries & via kernels;

Lecture 3 – Support Vector Machines 8 / 58

Expected Out-of-Sample Performance for 'Any Line'

- We believe there is a (linear) pattern to be detected
 - Assumption: linearly separable data (later non-separable cases)
 - Performance question: What is the optimal line (decision boundary)?
 - E.g. green data: $\left\{ \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \begin{pmatrix} 3 \\ 1 \end{pmatrix}, \begin{pmatrix} 4 \\ 1 \end{pmatrix} \right\}$ red data: $\left\{ \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right\}$




(PLA gives us just any line as soon as all samples are correctly classified)

(How can we craft a margin expressing 'furthest away')

- Intuition tells us just 'furthest away' from the closest points is a good position for the line – why?

Lecture 3 – Support Vector Machines 9 / 58

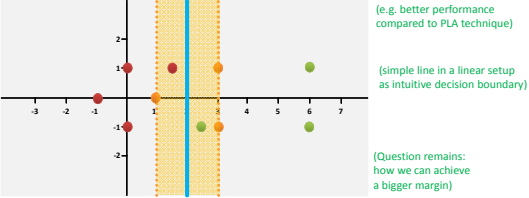
Exercises



Lecture 3 – Support Vector Machines 10 / 58

Expected Out-of-Sample Performance for ‘Best Line’

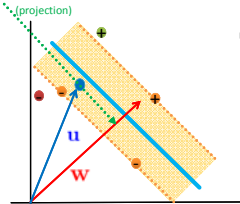
- The line with a ‘bigger margin’ seems to be better – but why?
 - Intuition: chance is higher that a new point will still be correctly classified
 - Fewer hypothesis possible: constrained by sized margin (cf. Lecture 2)
 - Idea: achieving good ‘out-of-sample’ performance is goal (cf. Lecture 2)



Lecture 3 – Support Vector Machines 11 / 58

Geometric SVM Interpretation and Setup (1)

- Think ‘simplified coordinate system’ and use ‘Linear Algebra’
 - Many other samples are removed (red and green not SVs) ● ⊕
 - Vector \mathbf{w} of ‘any length’ perpendicular to the decision boundary
 - Vector \mathbf{u} points to an unknown quantity (e.g. new sample to classify)
 - Is \mathbf{u} on the left or right side of the decision boundary?



- Dot product $\mathbf{w} \cdot \mathbf{u} \geq C; C = -b$
 - With \mathbf{u} takes the projection on the \mathbf{w}
 - Depending on where projection is it is left or right from the decision boundary
 - Simple transformation brings decision rule:
 - ① $\mathbf{w} \cdot \mathbf{u} + b \geq 0 \rightarrow$ means ●
 - (given that b and \mathbf{w} are unknown to us)
 - (constraints are not enough to fix particular b or \mathbf{w} , need more constraints to calculate b or \mathbf{w})

Lecture 3 – Support Vector Machines 12 / 58

Geometric SVM Interpretation and Setup (2)

- Creating our constraints to get b or w computed
 - First constraint set for positive samples \oplus $w \cdot x_+ + b \geq 1$
 - Second constraint set for negative samples \ominus $w \cdot x_- + b \leq -1$
 - For **mathematical convenience** introduce variables (i.e. labelled samples) $y_i = +$ for \oplus and $y_i = -$ for \ominus
- Multiply equations by y_i
 - Positive samples: $y_i(x_i \cdot w + b) \geq 1$
 - Negative samples: $y_i(x_i \cdot w + b) \geq 1$
 - Both **same** due to $y_i = +$ and $y_i = -$ (brings us mathematical convenience often quoted)

$y_i(x_i \cdot w + b) - 1 \geq 0$
 (additional constraints just for support vectors, itself helps)

② $y_i(x_i \cdot w + b) - 1 = 0$

Lecture 3 - Support Vector Machines 13 / 58

Geometric SVM Interpretation and Setup (3)

- Determine the 'width of the margin'
 - Difference between positive and negative SVs: $x_+ - x_-$
 - Projection of $x_+ - x_-$ onto the vector w
 - The vector w is a normal vector, magnitude is $\|w\|$

(Dot product of two vectors is a scalar, here the width of the margin)

- Unit vector is helpful for 'margin width'
 - Projection (dot product) for margin width: $(x_+ - x_-) \cdot \frac{w}{\|w\|}$ (unit vector)

$1 - b \quad 1 + b$

② $y_i(x_i \cdot w + b) - 1 = 0$

Lecture 3 - Support Vector Machines 14 / 58

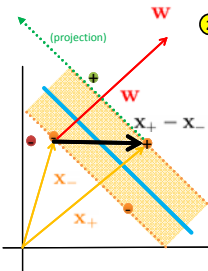
Constrained Optimization Steps SVM (1)

- Use 'constraint optimization' of mathematical toolkit
 - Idea is to 'maximize the width' of the margin: $\max \frac{2}{\|w\|}$ (drop the constant 2 is possible here)
 - $\Rightarrow \max \frac{1}{\|w\|}$ (equivalent)
 - $\Rightarrow \min \|w\|$ (equivalent for max)
 - $\Rightarrow \min \frac{1}{2} \|w\|^2$ (mathematical convenience) ③
- Next: Find the extreme values
 - Subject to constraints
 - ② $y_i(x_i \cdot w + b) - 1 = 0$

Lecture 3 - Support Vector Machines 15 / 58

Constrained Optimization Steps SVM (2)

- Use 'Lagrange Multipliers' of mathematical toolkit
 - Established tool in 'constrained optimization' to find function extremum
 - 'Get rid' of constraints by using Lagrange Multipliers ④



② $y_i(x_i \cdot w + b - 1) = 0$

- Introduce a multiplier for each constraint

$$\mathcal{L}(\alpha) = \frac{1}{2} \|w\|^2 - \sum \alpha_i [y_i(x_i \cdot w + b) - 1]$$

(interesting: non zero for support vectors, rest zero)
- Find derivatives for extremum & set 0
 - But two unknowns that might vary
 - First differentiate w.r.t. w
 - Second differentiate w.r.t. b

Lecture 3 - Support Vector Machines 16 / 58

Constrained Optimization Steps SVM (3)

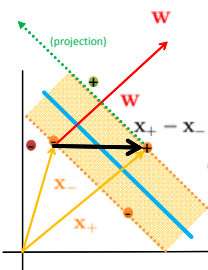
- Lagrange gives: $\mathcal{L}(\alpha) = \frac{1}{2} \|w\|^2 - \sum \alpha_i [y_i(x_i \cdot w + b) - 1]$
 - First differentiate w.r.t w

$$\frac{\partial \mathcal{L}}{\partial w} = w - \sum \alpha_i y_i x_i = 0$$
 - Simple transformation brings:

$$w = \sum \alpha_i y_i x_i$$

(i.e. vector is linear sum of samples)
(recall: non zero for support vectors, rest zero → even less samples)
 - Second differentiate w.r.t. b

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum \alpha_i y_i = 0 \Rightarrow \sum \alpha_i y_i = 0$$



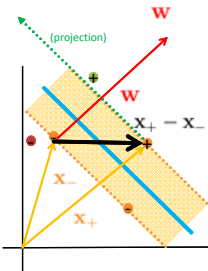
Lecture 3 - Support Vector Machines 17 / 58

Constrained Optimization Steps SVM (4)

- Lagrange gives: $\mathcal{L}(\alpha) = \frac{1}{2} \|w\|^2 - \sum \alpha_i [y_i(x_i \cdot w + b) - 1]$
 - Find minimum
 - Quadratic optimization problem
 - Take advantage of ⑤ $w = \sum \alpha_i y_i x_i$

$$\mathcal{L} = \frac{1}{2} (\sum \alpha_i y_i x_i) \cdot (\sum \alpha_j y_j x_j) - \sum \alpha_i y_i x_i \cdot (\sum \alpha_j y_j x_j) - \sum \alpha_i y_i b + \sum \alpha_i$$

(b constant in front sum) ⑤ $\sum \alpha_i y_i = 0$



Lecture 3 - Support Vector Machines 18 / 58

Constrained Optimization Steps SVM (5)

- Rewrite formula: $\mathcal{L} = \frac{1}{2} \sum \alpha_i y_i x_i \cdot (\sum \alpha_j y_j x_j)$ (the same)

 $-\sum \alpha_i y_i x_i \cdot (\sum \alpha_j y_j x_j)$

 $-\sum \alpha_i y_i b + \sum \alpha_i$ (was 0)

 (results in)

 $\mathcal{L} = \sum \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i \cdot x_j$ (optimization depends only on dot product of samples)

 Equation to be solved by some quadratic programming package

Lecture 3 - Support Vector Machines 19 / 58

Use of SVM Classifier to Perform Classification

- Use findings for decision rule

 $w = \sum \alpha_i y_i x_i$

 $w \cdot u + b \geq 0$

 $\sum \alpha_i y_i x_i \cdot u_i + b \geq 0$ (decision rule also depends on dot product)

Lecture 3 - Support Vector Machines 20 / 58

Exercises

Lecture 3 - Support Vector Machines 21 / 58

LibSVM – Defacto Standard SVM Implementation

- Free available tool
 - Includes Sequential Minimal Optimization (SMO) implementation

[2] LibSVM Webpage
Lecture 3 – Support Vector Machines 22 / 58



LibSVM – Download

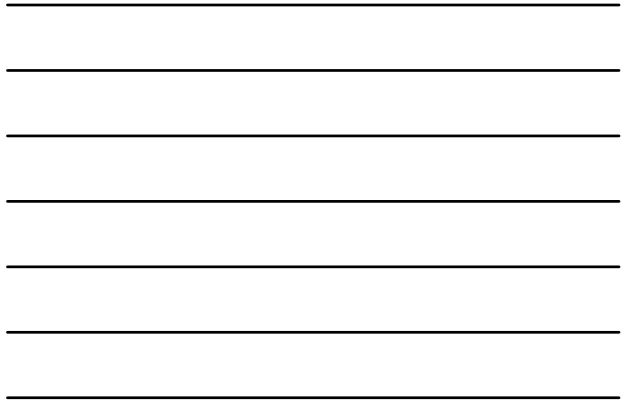
- Download tar.gz (or in Windows zip bundle)

[2] LibSVM Webpage

- Put package in a folder of your choice
 - Alternatively copy file to your usual working environment

```
adminuser@linux-8djq:~/tools> scp libsvm-3.21.tar.gz mriedel@ureca.fz-juelich.de:/home/zam/mriedel
libsvm-3.21.tar.gz 100% 827KB 827.4KB/s 00:00

-bash-4.2$ ls -al
total 64
drwxr-xr-x 2 mriedel zam 512 Jul 6 20:00 .
drwxr-xr-x 29 mriedel zam 32768 Jul 6 19:58 ..
-rw-r--r-- 1 mriedel zam 847291 Jul 6 20:00 libsvm-3.21.tar.gz
-bash-4.2$ pwd
/home/zam/mriedel/serialtools
Lecture 3 – Support Vector Machines 23 / 58
```



LibSVM – Unpack the Bundle

- Untar (or Unzip in Windows)

```
/home/zam/mriedel/serialtools
-bash-4.2$ tar -zxvf libsvm-3.21.tar.gz
libsvm-3.21/
libsvm-3.21/COPYRIGHT
libsvm-3.21/svm.cpp
libsvm-3.21/README
libsvm-3.21/Makfile.win
libsvm-3.21/svm.h
libsvm-3.21/svm.h
libsvm-3.21/mkafila
libsvm-3.21/java/
libsvm-3.21/java/ava_too.java
libsvm-3.21/java/ava_scale.java
libsvm-3.21/java/libsvm
libsvm-3.21/java/libsvm/svm_model.java
libsvm-3.21/java/libsvm/svm.h
libsvm-3.21/java/libsvm/svm_problem.java
libsvm-3.21/java/libsvm/svm.java
libsvm-3.21/java/libsvm/svm_node.java
libsvm-3.21/java/libsvm/svm_parameter.java
libsvm-3.21/java/libsvm/svm_print_interface.java
libsvm-3.21/java/ava_train.libsvm
libsvm-3.21/java/Makfile
libsvm-3.21/java/test_svm1.html
libsvm-3.21/java/ava_predict.java
libsvm-3.21/Makfile
libsvm-3.21/windows/
libsvm-3.21/windows/svm-tyt.exe
libsvm-3.21/windows/svm-scale.exe
libsvm-3.21/windows/ava_train.makfile
libsvm-3.21/windows/libsvmwrite.makfile
libsvm-3.21/windows/libsvm.vcl
/home/zam/mriedel/serialtools/libsvm-3.21
-bash-4.2$ ls -al
total 548
drwxr-xr-x 8 mriedel zam 32768 Dec 14 2015 .
drwxr-xr-x 3 mriedel zam 512 Jul 6 20:03 ..
-rw-r--r-- 1 mriedel zam 1407 Dec 14 2015 COPYRIGHT
-rw-r--r-- 1 mriedel zam 83089 Dec 14 2015 FAQ.html
-rw-r--r-- 1 mriedel zam 23570 Dec 14 2015 mkafila
drwxr-xr-x 3 mriedel zam 512 Dec 14 2015 java
-rw-r--r-- 2 mriedel zam 752 Dec 14 2015 Makfile
-rw-r--r-- 1 mriedel zam 1196 Dec 14 2015 Makfile.win
drwxr-xr-x 2 mriedel zam 512 Dec 14 2015 svm.h
drwxr-xr-x 2 mriedel zam 512 Dec 14 2015 svm.h
-rw-r--r-- 1 mriedel zam 28079 Dec 14 2015 README
-rw-r--r-- 1 mriedel zam 6496 Dec 14 2015 svm.cpp
-rw-r--r-- 1 mriedel zam 477 Dec 14 2015 svm.def
-rw-r--r-- 1 mriedel zam 3382 Dec 14 2015 svm.h
-rw-r--r-- 1 mriedel zam 5536 Dec 14 2015 svm-predict.c
-rw-r--r-- 1 mriedel zam 8096 Dec 14 2015 svm-scale.c
drwxr-xr-x 5 mriedel zam 512 Dec 14 2015 svm-tyt
-rw-r--r-- 1 mriedel zam 8096 Dec 14 2015 svm-train.c
drwxr-xr-x 2 mriedel zam 512 Dec 14 2015 tools
drwxr-xr-x 2 mriedel zam 512 Dec 14 2015 windows
Lecture 3 – Support Vector Machines 24 / 58
```



LibSVM – Make (only in UNIX)

- Use make to generate executables (needs g++ compiler)

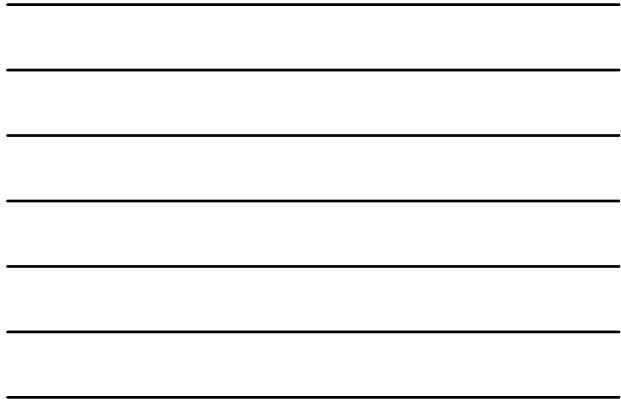
```
-bash-4.2$ pwd
/home/zam/riedel/serialtools/libsvm-3.21
-bash-4.2$ make
g++ -Wall -Mconversion -O3 -fPIC -c svm.cpp
g++ -Wall -Mconversion -O3 -fPIC svm-train.c svm.o -o svm-train -lm
g++ -Wall -Mconversion -O3 -fPIC svm-predict.c svm.o -o svm-predict -lm
g++ -Wall -Mconversion -O3 -fPIC svm-scale.c -o svm-scale
```

- Check executables important for us

```
-bash-4.2$ pwd
/home/zam/riedel/serialtools/libsvm-3.21
-bash-4.2$ ls -al
total 808
drwxr-xr-x 8 erinabl sam 32768 Jul 6 20:05 .
drwxr-xr-x 2 erinabl sam 512 Jul 6 20:02 ..
-rw-r--r-- 1 erinabl sam 1487 Dec 14 2015 COPYING
-rw-r--r-- 1 erinabl sam 40980 Dec 14 2015 FAQ.html
-rw-r--r-- 1 erinabl sam 27678 Dec 14 2015 heart_scale
drwxr-xr-x 2 erinabl sam 512 Dec 14 2015 java
-rw-r--r-- 1 erinabl sam 712 Dec 14 2015 Makefile
drwxr-xr-x 2 erinabl sam 1150 Dec 14 2015 MatlabCode
drwxr-xr-x 2 erinabl sam 512 Dec 14 2015 matlab
drwxr-xr-x 2 erinabl sam 512 Dec 14 2015 python
-rw-r--r-- 1 erinabl sam 20679 Dec 14 2015 README
-rw-r--r-- 1 erinabl sam 14838 Dec 14 2015 svm.cpp
-rw-r--r-- 1 erinabl sam 477 Dec 14 2015 svm.daf
-rw-r--r-- 1 erinabl sam 3362 Dec 14 2015 svm.h
-rw-r--r-- 1 erinabl sam 78275 Jul 6 20:05 svm-train.o (use in testing phase)
-rw-r--r-- 1 erinabl sam 18267 Jul 6 20:06 svm-train.o.gch
-rw-r--r-- 1 erinabl sam 8538 Dec 14 2015 svm-scale.o
-rw-r--r-- 1 erinabl sam 78959 Jul 6 20:06 svm-scale.o (use in training phase)
drwxr-xr-x 2 erinabl sam 512 Dec 14 2015 tools
drwxr-xr-x 2 erinabl sam 512 Dec 14 2015 windows
```

[2] LibSVM Webpage
Lecture 3 – Support Vector Machines

25 / 58



LibSVM – svm-train Parameters

- Important parameters (training phase)

```
-bash-4.2$ ./svm-train
Usage: svm-train [options] training_set_file [model_file]
-s svm_type : set type of SVM (default 0)
  0 -- C-SVC (multi-class classification)
  1 -- nu-SVC (multi-class classification)
  2 -- one-class SVM (regression)
  3 -- epsilon-SVR (regression)
-t kernel_type : set type of kernel function (default 2)
  0 -- linear: u^T v
  1 -- polynomial: (gamma*u^T v + coef0)^degree
  2 -- radial basis function: exp(-gamma*|u-v|^2)
  3 -- sigmoid: tanh(gamma*u^T v + coef0)
  4 -- precomputed kernel (kernel values in training_set_file)
-d degree : set degree in kernel function (default 3)
-g gamma : set gamma in kernel function (default 1/num_features)
-r coef0 : set coef0 in kernel function (default 0)
-c cost : set the parameter C of C-SVC, epsilon-SVM, and nu-SVM (default 1) (Regularization Parameter)
-h nr : set the parameter nu of nu-SVC, epsilon-SVM, and nu-SVM (default 0.5)
-p epsilon : set the epsilon in loss function of epsilon-SVR (default 0.1)
-m cachesize : set cache memory size in MB (default 100)
-e epsilon : set tolerance of termination criterion (default 0.001)
-h shrinking : whether to use the shrinking heuristics, 0 or 1 (default 1)
-b probability_estimates : whether to train a SVC or SVR model for probability estimates, 0 or 1 (default 1)
-wd weight : set the parameter C of class i to weight^C, for C-SVC (default 1)
-v n : n-fold cross validation mode
-q : quiet mode (no outputs)
```

Training Examples
 $(X_1, Y_1), \dots, (X_N, Y_N)$

[2] LibSVM Webpage
Lecture 3 – Support Vector Machines

26 / 58



LibSVM – svm-predict Parameters

- Important parameters (testing phase)

```
-bash-4.2$ ./svm-predict
Usage: svm-predict [options] test_file model_file output_file
-p probability_estimates : whether to predict probability estimates, 0 or 1 (default 0); for one-class SVM only 0 is supported
-q : quiet mode (no outputs)
```

(the model file is generated in the training phase → the support vectors found in optimization)

(test file is a testing dataset set aside to be used once training is finished)

(output file gives us indications how each sample was classified)

Testing Examples
 $(X_1, Y_1), \dots, (X_N, Y_N)$

Lecture 3 – Support Vector Machines

27 / 58

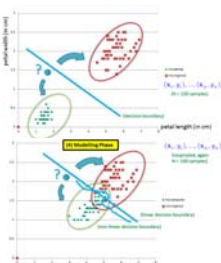


Data Preparation Phase

- Copy IRIS Dataset in your working environment

```
adminuser@linux-8dgi:~/data> scp iris.scale.mriedel@jurca.fz.juelich.de:/home/zam/mriedel/iris.scale
~/home/zam/mriedel/datasets
-bash-4.2$ ls -al
total 64
drwxr-xr-x 2 mriedel zam 512 Jul 6 21:53
drwxr-xr-x 8 mriedel zam 32768 Jul 6 21:53
-rw-r--r-- 1 mriedel zam 4564 Jul 6 21:53 iris.scale
```

- Dataset Two-class problem, linearly seperable
 - Dataset Iris Setosa (class 1) and Iris Virginica (class 3)
 - iris-class1and3-training(20)/testing(30)
- Dataset Two-class problem, not linearly seperable
 - Dataset Iris Vericolor (class 2) and Iris Virginica (class3)
 - iris-class2and3-training(20)/testing(30)



Lecture 3 - Support Vector Machines

31 / 58



Data Download

- Download the following files from B2Share



Iris Dataset LibSVM Format Preprocessing

```
03 July 2016
http://www.b2share.eu/

Abstract: LIBSVM Machine Learning Repository
URL: Dataset
Website: original and iris.scale
3 classes, all samples each class
iris.scale(iris)
-only locally accessible data
-class 1 and 3 sampling
-copy classes
-iris.scale(iris) training/testing
-iris.scale(iris) testing
-copy class 2
-iris.scale(iris) training/testing
-iris.scale(iris) testing
-copy class 3 only

Keywords: LIBSVM(iris), Flowers(iris)
The second appears in other collections:
class=2

-bash-4.2$ pwd
~/home/zam/mriedel/datasets
-bash-4.2$ ls -al
total 92
drwxr-xr-x 2 mriedel zam 512 Jul 6 22:16
drwxr-xr-x 8 mriedel zam 32768 Jul 6 22:16
-rw-r--r-- 1 mriedel zam 4562 Jul 6 21:58 iris-class1and3
-rw-r--r-- 1 mriedel zam 4736 Jul 6 22:16 iris-class1and3-training
-rw-r--r-- 1 mriedel zam 3892 Jul 6 22:16 iris-class1and3-testing
-rw-r--r-- 1 mriedel zam 4058 Jul 6 22:16 iris-class2and3-training
-rw-r--r-- 1 mriedel zam 3981 Jul 6 22:16 iris-class2and3-testing
-rw-r--r-- 1 mriedel zam 3188 Jul 6 22:16 iris-class2and3-testing
-rw-r--r-- 1 mriedel zam 4954 Jul 6 21:54 iris.scale.original
```

Lecture 3 - Support Vector Machines

32 / 58



Training Phase: linearly seperable case (iris-class1and3)

- Use svm-train (c<=0 not allowed)

```
-bash-4.2$ more svm-train1-3.sh
./svm-train -t 0 -c 1 /home/zam/mriedel/datasets/iris-class1and3-training
-bash-4.2$ ./svm-train1-3.sh
#
optimization finished, #iter = 11
nu = 0.935490
obj = -0.769742, rho = 0.447384
nSV = 4, nBSV = 0
Total nSV = 4

-bash-4.2$ ls -al
total: 896
drwxr-xr-x 8 mriedel zam 32768 Jul 6 22:25
drwxr-xr-x 2 mriedel zam 512 Jul 6 20:03
-rw-r--r-- 1 mriedel zam 1487 Dec 14 2015 COPYINGRIGHT
-rw-r--r-- 1 mriedel zam 6369 Dec 14 2015 README
-rw-r--r-- 1 mriedel zam 354 Jul 6 22:25 iris-class1and3-training.model
```

- Check model file

```
-bash-4.2$ more iris-class1and3-training.model
svm_model
svm_type C_SVC
nr_class 3
nr_inst 150
svm_svm_type C
svm_kernel rbf
gamma 0.1
nu 0.935490
c 1
svm_coef 0.000000000000 0.000000000000 0.000000000000
svm_weight 0.000000000000 0.000000000000 0.000000000000
svm_bias 0.000000000000 0.000000000000 0.000000000000
svm_gamma 0.100000000000 0.100000000000 0.100000000000
svm_coef 0.000000000000 0.000000000000 0.000000000000
svm_weight 0.000000000000 0.000000000000 0.000000000000
svm_bias 0.000000000000 0.000000000000 0.000000000000
svm_gamma 0.100000000000 0.100000000000 0.100000000000
svm_coef 0.000000000000 0.000000000000 0.000000000000
svm_weight 0.000000000000 0.000000000000 0.000000000000
svm_bias 0.000000000000 0.000000000000 0.000000000000
svm_gamma 0.100000000000 0.100000000000 0.100000000000
```

Lecture 3 - Support Vector Machines

33 / 58



Maximal Margin Classifier – Margin Performance

- Classification technique
 - Classify testset samples based on which side of the maximal margin hyperplane they are lying

sign of $f(x^*) = \beta_0 + \beta_1 x_1^* + \beta_2 x_2^* + \dots + \beta_p x_p^*$
 (β s are the coefficients of the maximal margin hyperplane)

- Assuming that a classifier that has a large margin on the training data will also have a large margin on the test data (cf. also 'the intuitive notion')
- Testset samples will be thus correctly classified

modified from [1] An Introduction to Statistical Learning
 Lecture 3 – Support Vector Machines 37 / 58

Maximal Margin Classifier – Support Vector Term

- Observation
 - Three data points lie on the edge of margin (somewhat special data points)
 - Dashed lines indicating the width of the margin (very interesting to know)
 - Margin width is the distance from the special data points to the hyperplane (hyperplane depends directly on small data subset; SV points)

modified from [1] An Introduction to Statistical Learning

- Points that lie on the edge of the margin are named support vectors (SVs) in p-dimensional space
- SVs 'support' the maximal margin hyperplane: if SVs are moved → the hyperplane moves as well

Lecture 3 – Support Vector Machines 38 / 58

Maximal Margin Classifier – Optimization and W Vector

- Which weight w maximizes the margin?
 - Margin is just a distance from 'a line to a point', goal is to minimize w
 - Pick x_n as the nearest data point to the line (or hyper-plane)...

(distance between two dashed planes is $2 / ||w||$)

- Reduce the problem to a 'constraint optimization problem' (vector w are the β coefficients)

$$\text{maximize } M$$

$$\text{subject to } \sum_{j=1}^p \beta_j^2 = 1$$

(for points on plane w must be 0, interpret k as length of w)

$$\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} = 0 \quad k(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) = 0 \text{ for any } k \neq 0$$

- Support vectors achieve the margin and are positioned exactly on the boundary of the margin

Lecture 3 – Support Vector Machines 39 / 58

Maximal Margin Classifier – Solving and Limitations

- Solving constraint optimization problem chooses coefficients that maximize M & gives hyperplane
- Solving this problem efficiently is possible techniques like sequential minimal optimization (SMO)
- Maximal margin classifiers use a hard-margin & thus only work with exact linearly separable data

modified from [1] An Introduction to Statistical Learning

- Limitation
 - Non linearly separable data (given mostly in practice)
 - Optimization problem has no solution $M > 0$ (think point moves over plane)
 - No separating hyperplane can be created (classifier can not be used)

(no error allowed, a 'hard margin')

(allow some error the margin will be bigger, ... maybe better E_{test})

(move effects the margin)

(no exact separation possible)

(... but with allowing some error maybe, a 'soft margin'...)

Lecture 3 – Support Vector Machines 43 / 58

Support Vector Classifier

Lecture 3 – Support Vector Machines 44 / 58

Support Vector Classifiers – Motivation

- Support Vector Classifiers develop hyperplanes with soft-margins to achieve better performance
- Support Vector Classifiers aim to avoid being sensitive to individual observations (e.g. outliers)

[1] An Introduction to Statistical Learning

- Approach
 - Generalization of the 'maximal margin classifier' (get most & but not all training data correctly classified)
 - Include non-separable cases with a soft-margin (almost instead of exact)
 - Being more robust w.r.t. extreme values (e.g. outliers) allowing small errors

(outlier)

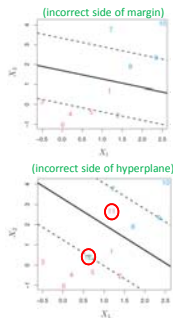
(significant reduction of the maximal margin)

(overfitting, cf Lecture 10: maximal margin classifier & hyperplane is very sensitive to a change of a single data point)

Lecture 3 – Support Vector Machines 45 / 58

Support Vector Classifiers – Modified Technique Required

- Previous classification technique reviewed
 - Seeking the largest possible margin, so that every data point is...
 - ... on the **correct side of the hyperplane**
 - ... on the **correct side of the margin**
- Modified classification technique
 - Enable 'violations of the hard-margin' to achieve 'soft-margin classifier'
 - Allow violations means: allow some data points to be...
 - ... on the **incorrect side of the margin**
 - ... even on the **incorrect side of the hyperplane** (which leads to a **misclassification of that point**)



[1] An Introduction to Statistical Learning (SVMs refined: data points that lie directly on the margin or on the wrong side of the margin for their class are the SVs → affect support vector classifier)
 Lecture 3 – Support Vector Machines 46 / 58

Support Vector Classifier – Optimization Problem Refined

- Optimization Problem
 - Still **maximize margin M**
 - Refining constraints to include **violation of margins**
 - Adding **slack variables** $\epsilon_1, \dots, \epsilon_n$ (allow datapoints to be on the wrong side of the margin or hyperplane)
- $$\text{maximize } M$$
- $$\text{subject to } \sum_{j=1}^p \beta_j^2 = 1,$$
- $$y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M(1 - \epsilon_i) \text{ (slightly violate the Margin)}$$
- $$\epsilon_i \geq 0, \sum_{i=1}^n \epsilon_i \leq C$$
- (C is used here to bound the error)
- C Parameter & slacks**
 - C bounds the sum of the slack variables $\epsilon_1, \dots, \epsilon_n$
- [1] An Introduction to Statistical Learning

C determines the number & severity of violations that will be tolerated to margin and hyperplane
 Interpret C as budget (or costs) for the amount that the margin can be violated by n data points
 Lecture 4 provides details on validation methods that provides a value for C in a principled way

Lecture 3 – Support Vector Machines 47 / 58

Support Vector Classifier – Understanding the Slack

- Allowing some errors or **violations of the margin**
-
- (As budget C increases, the classifier becomes more tolerant of violations of the margin, margin will be thus wider)
- [4] SVMs
- (These all correct samples don't matter contributing to a robust classifier and unique property of SVCs)

Lecture 3 – Support Vector Machines 48 / 58

Exercises



Training Phase: non-linearly seperable case (iris-class2and3)

- Use svm-train (c<=0 not allowed, change value, what happens?)

```

-bash-4.2$ more svm-train2-3.sh
./svm-train -t 0 -c 1 /home/zam/mriedel/datasets/iris-class2and3-training

-bash-4.2$ ./svm-train2-3.sh
svm-4.28 ts -t
total: 800
dfrw@fira:1 svm@k81 sam 32768 Jul 6 22:36
dfrw@fira:2 svm@k81 sam 512 Jul 6 22:03
dfrw@fira:3 svm@k81 sam 1637 Dec 14 2015 CPWRIGet
dfrw@fira:4 svm@k81 sam 8088 Dec 14 2015 FQJmpet
dfrw@fira:5 svm@k81 sam 27678 Dec 14 2015 heart_scale
dfrw@fira:6 svm@k81 sam 7884 Jul 6 22:06 iris-class2and3-training.model
dfrw@fira:7 svm@k81 sam 1490 Jul 6 22:06 iris-class2and3-training.model
dfrw@fira:8 svm@k81 sam 722 Dec 14 2015 Manafila
dfrw@fira:9 svm@k81 sam 1190 Dec 14 2015 Manafila.pdf
dfrw@fira:10 svm@k81 sam 512 Dec 14 2015 mailub
dfrw@fira:11 svm@k81 sam 512 Dec 14 2015 sifham
dfrw@fira:12 svm@k81 sam 28678 Dec 14 2015 REQONE
dfrw@fira:13 svm@k81 sam 120 Jul 6 22:32 ResNetA.txt
dfrw@fira:14 svm@k81 sam 68678 Dec 14 2015 svm_0ap
dfrw@fira:15 svm@k81 sam 477 Dec 14 2015 svm_0AP
dfrw@fira:16 svm@k81 sam 3962 Dec 14 2015 svm_0h
dfrw@fira:17 svm@k81 sam 300274 Jul 6 20:05 svm_0
dfrw@fira:18 svm@k81 sam 78278 Jul 6 20:05 svm_0redist
dfrw@fira:19 svm@k81 sam 419 Jul 6 22:31 svm_0redist-3.0h
dfrw@fira:20 svm@k81 sam 118 Jul 6 22:35 svm_0redist-3.0h
dfrw@fira:21 svm@k81 sam 78278 Jul 6 20:05 svm_0redist-3.0h
dfrw@fira:22 svm@k81 sam 18567 Jul 6 20:05 svm_0redist-3.0h
dfrw@fira:23 svm@k81 sam 8538 Dec 14 2015 svm_0redist-3.0h
dfrw@fira:24 svm@k81 sam 512 Dec 14 2015 svm_0redist-3.0h
dfrw@fira:25 svm@k81 sam 78278 Jul 6 20:05 svm_0redist-3.0h
dfrw@fira:26 svm@k81 sam 76 Jul 6 22:24 svm_train2-3.0h
dfrw@fira:27 svm@k81 sam 76 Jul 6 22:35 svm_train2-3.0h
dfrw@fira:28 svm@k81 sam 8088 Dec 14 2015 svm_train2.c
dfrw@fira:29 svm@k81 sam 512 Dec 14 2015 svm_train2.c
dfrw@fira:30 svm@k81 sam 512 Dec 14 2015 svm_train2.c

```

- Check model file
 - Next page, because many support vectors!

Model File: non-linearly seperable case (iris-class2and3)

- Many SVs / sample
 - (careful – indicator for problems)
 - In the linear case we know from looking at the data it still be ok in this case
 - Linear SVM worked: PLA instead would not be able to stop with this dataset

```

svm-4.28 more iris-class2and3-training.model
svm_7806_c_svm
kernel_type: Linear
nr_classes: 2
nr_sv: 27
rho: -1.69284
margin: 2.9
nr_sv: 14 13
SV:
1 10.5 3.0 254237 4.0 0.003333
1 10.166667 3.0 186441 4.0 1.06667
1 10.44444 2.0 0.003334 3.0 322934 4.0 1.06667
1 10.33333 2.0 76.0 0.003333 4.0 4.02576e-08
1 10.32222 2.0 0.33333 3.0 320399 4.0 1.06667
1 10.22222 2.0 0.33333 3.0 186441 4.0 4.02576e-08
1 10.11111 2.0 0.003333 3.0 254237 4.0 2.5
1 10.27778 2.0 2.25 3.0 220326 4.0 4.02576e-08
0.300857881625212 1.0 0.5 2.0 0.416667 2.0 0.016667 4.0 0.003333
1 10.11111 2.0 1.06667 3.0 186441 4.0 1.06667
1 1.1 3.04545 0.7 2.0 2.25 3.0 254237 4.0 0.003333
0.3223908917076 1.0 0.33333 3.0 0.003334 3.0 177342 4.0 0.003333
1 10.27778 2.0 1.06667 3.0 186441 4.0 1.06667
1 10.00000 2.0 0.33333 3.0 186441 4.0 1.06667
1 10.66667 2.0 1.06667 3.0 186441 4.0 1.06667
0.0420374010240 1.0 2.22222 2.0 0.00000 4.0 0.00333
1 10.20222 2.0 1.06667 3.0 325424 4.0 4.16667
1 10.00000 2.0 0.33333 3.0 305302 4.0 1.06667
1 10.11111 2.0 0.41667 3.0 322934 4.0 4.16667
1 10.00000 2.0 0.33333 3.0 305302 4.0 1.06667
1 1.1 3.04545 0.7 2.0 1.06667 3.0 322934 4.0 4.16667
1 10.11111 2.0 1.06667 3.0 186441 4.0 1.06667
1 1.1 3.04545 0.7 2.0 0.5 3.0 305302 4.0 0.003333
1 10.16667 2.0 0.003334 3.0 325424 4.0 4.16667
1 10.00000 2.0 0.16667 3.0 305302 4.0 1.06667
1 10.11111 2.0 1.06667 3.0 186441 4.0 1.06667

```

- Generalization measure: #SVs as 'in-sample quantity' → 10SVs/1000 samples ok, 500SVs/1000 bad
- Reasoning towards overfitting due to a large number of SVs (fit many, small margin, gives bad E_{out})
