





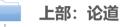
## 既是世间法、自当有分别

艾新波 / 2018 • 北京



#### 课程体系







- 第2章 所谓学习、归类而已
- 第3章 格言联璧话学习
- 第4章 源于数学、归于工程
- 中部:执具
  - 第5章 工欲善其事必先利其器
  - 第6章 基础编程
  - 第7章 数据对象

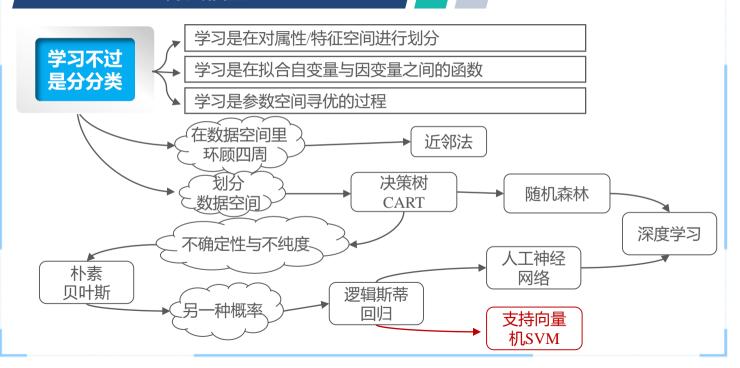




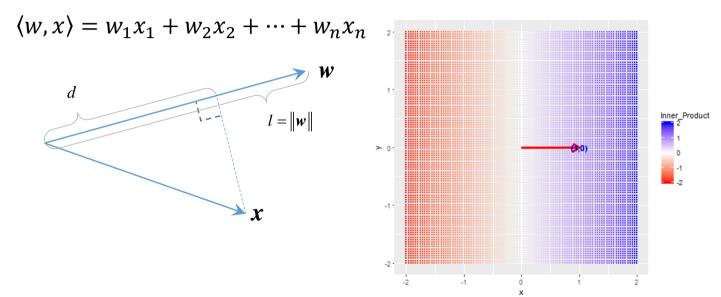


- 第10章 观数以形
- 第11章 相随相伴、谓之关联
  - 🗐 第12章 既是世间法、自当有分别
  - 第13章 方以类聚、物以群分
  - 第14章 庐山烟雨浙江潮

## 算法模型



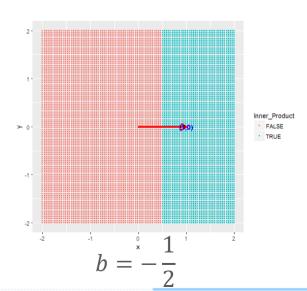
## 从内积说起

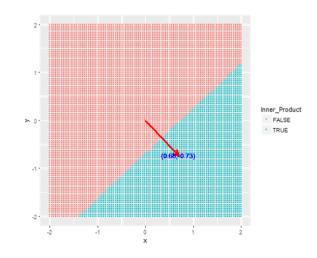


通俗讲: w与x的内积, 就是x在w上的投影长度, 乘以w的长度

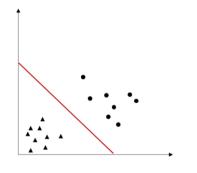
## 从内积到分类超平面

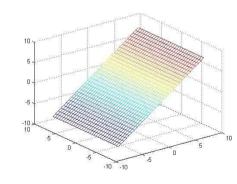
$$w^T x + b = 0$$
作为一个分类超平面:

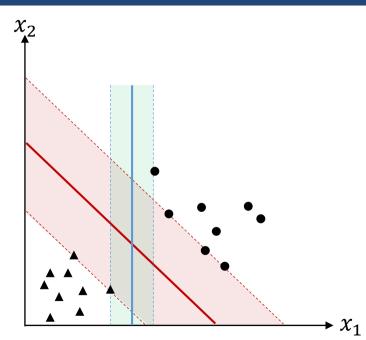




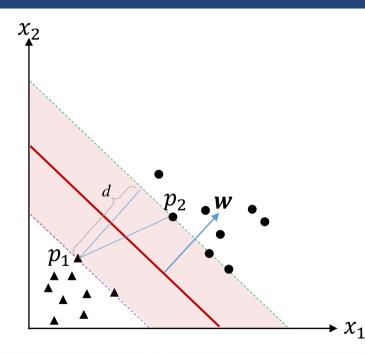
$$w^T x + b = 0$$







- 要将作图所示的两类点分开, 可以有无穷多个超平面
- 相比较而言,红色的部分有更 宽的边缘
- 口 "道路" 更宽、泛化能力更强



## 两个平行的边界为:

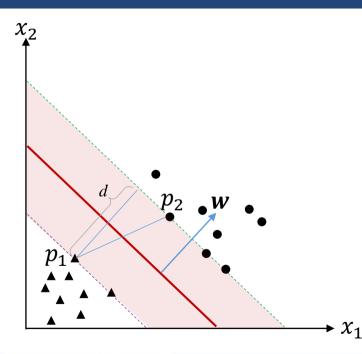
$$\mathbf{w}^T \mathbf{x} + b = 1$$
$$\mathbf{w}^T \mathbf{x} + b = -1$$

## 于是有:

$$\mathbf{w}^{T} (\mathbf{p}_{2} - \mathbf{p}_{1}) = 2$$

$$\Rightarrow \|\mathbf{w}\| \times d = 2$$

$$\Rightarrow d = \frac{2}{\|\mathbf{w}\|}$$



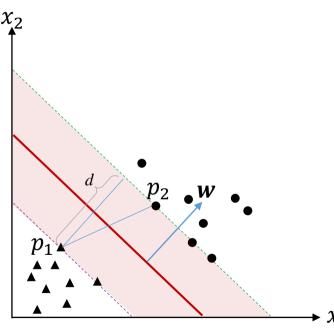
## 实现间隔最大化,等价于最小化 目标函数:

$$\min_{\mathbf{w}} \frac{\|\mathbf{w}\|^2}{2}$$

## 需要满足以下约束:

$$\mathbf{w}^{T} \mathbf{x} + b \ge 1, y_{i} = 1$$
  
 $\mathbf{w}^{T} \mathbf{x} + b \le -1, y_{i} = -1$ 

**即:** 
$$y_i(\mathbf{w}^T\mathbf{x}+b) \geq 1$$



## 通过拉格朗日法求解

## 对每个样本引入一个拉格朗日系数:

$$\alpha_i \geq 0, i = 1, 2, \dots, N$$

权向量最优解:  $\mathbf{w}^* = \sum_{i=1}^{N} \alpha_i^* y_i \mathbf{x}^{(i)}$ 

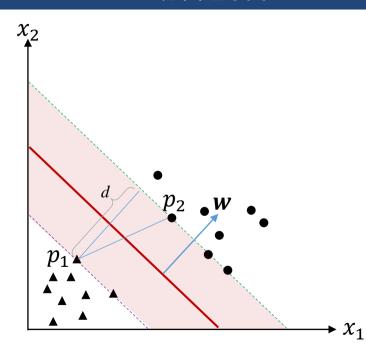
$$\Rightarrow \mathbf{w}^T \mathbf{x}^{(j)} + b$$

$$= \left(\sum_{i=1}^{N} \alpha_i^* y_i \boldsymbol{x}^{(i)}\right)^T \boldsymbol{x}^{(j)} + b$$

$$=\sum_{i=1}^{N}\alpha_{i}^{*}y_{i}\left(\boldsymbol{x}^{(i)}\right)^{T}\boldsymbol{x}^{(j)}+b$$

$$=\sum_{i=1}^{N}\alpha_{i}^{*}y_{i}\left\langle \boldsymbol{x}^{(i)},\boldsymbol{x}^{(j)}\right\rangle+b$$

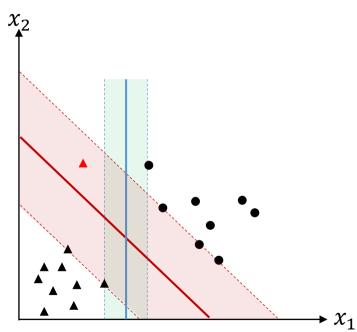
 $x_1$ 



#### 从决策函数可以看出:

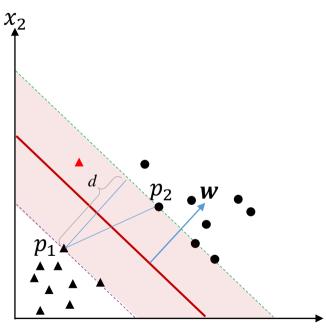
- 口 对于测试集的预测,只需要计算它与训练数据点的内积即可
- 一 准确地讲,是只需要计算它和支持向量的内积即可,因为非支持向量的α均为零

## 噪声与软间隔



- 口 红色的超平面,错分了一个样本,蓝色的超平面完全正确
- 口 但是,红色道路依然比蓝色的道路更可取
- 软间隔soft margin学习允许一 定训练错误的决策边界

## 噪声与软间隔



## 引入松弛变量:

$$\mathbf{w}^{T} \mathbf{x}^{(i)} + b = 1 - \xi_{i}, y_{i} = 1$$
  
 $\mathbf{w}^{T} \mathbf{x}^{(i)} + b = -1 + \xi_{i}, y_{i} = -1$ 

$$m{w} = m{x} + m{v} - m{1} + m{\varsigma}_i, \, m{y}_i - m{y}_i$$

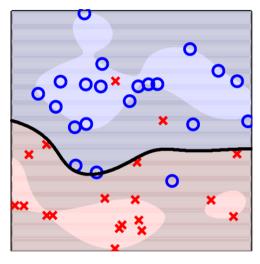
## 修正目标函数:

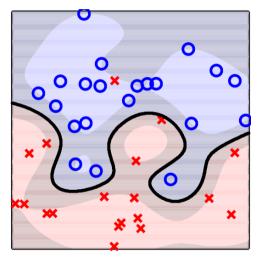
$$\frac{\|\mathbf{w}\|^2}{2} + C\left(\sum_{i=1}^N \xi_i\right)^k$$

C越大,对于噪声的容忍就越小

 $x_1$ 

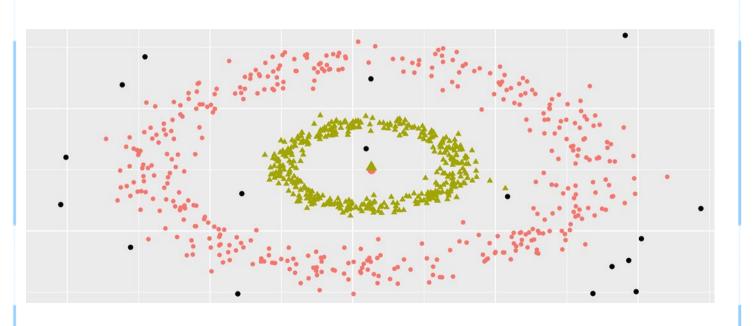
## 噪声与软间隔



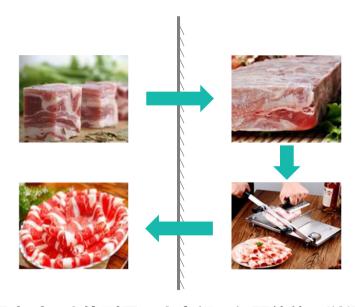


C=1 C=10

## 非线性可分问题

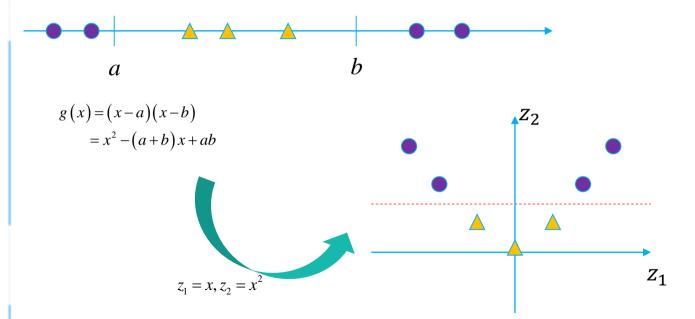


## 非线性可分:映射与变换



当问题不好解决时, 变换到另一个空间, 问题往往可以得到解决

## 支持向量机: 线性不可分



#### 支持向量机:线性不可分

对于非线性问题,需要做到:

第一,将属性空间映射到特征空间;

第二,在特征空间中采用线性分类器进行分类。

决策规则为: 
$$f(x) = \sum_{i=1}^{l} \alpha_i y_i \langle \phi(x^{(i)}), \phi(x) \rangle + b$$

其中,  $\phi$  为属性空间到特征空间的映射:  $\phi: \mathcal{X} \to \mathcal{F}$ 

#### 引入核方法:

一对向量的核函数等于变换后空间中计算这对向量的内积:  $K(x,z) = \langle \phi(x), \phi(z) \rangle$ 

也就是说,有了核函数,我们甚至不需要知道如何具体映射到高维空间

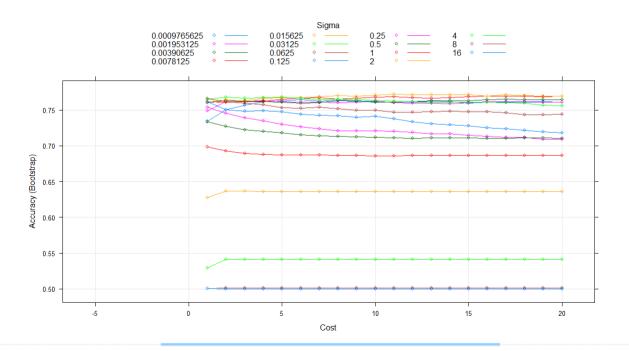
## 支持向量机:来自libsvm的一些建议

## Tips on practical use

- Note that SVMs may be very sensible to the proper choice of parameters, so allways check a range of parameter combinations, at least on a reasonable subset of your data.
- For classification tasks, you will most likely use C-classification with the RBF kernel (default), because of its good general performance and the few number of parameters (only two: C and  $\gamma$ ). The authors of libsvm suggest to try small and large values for C—like 1 to 1000—first, then to decide which are better for the data by cross validation, and finally to try several  $\gamma$ 's for the better C's.

```
library(kernlab)
imodel <- ksvm(wlfk~., data = cjb[train set idx, ])</pre>
predicted train <- predict(</pre>
  imodel, newdata = cjb[train set idx,], type = "response")
Metrics::ce(cjb$wlfk[train set idx], predicted train)
#> [1] 0.1497227
predicted test <- predict(imodel,</pre>
  newdata = cjb[-train set idx,], type = "response")
Metrics::ce(cjb$wlfk[-train set idx], predicted test)
#> [1] 0.1802575
```

```
library(caret)
svm grid <- expand.grid(</pre>
  sigma = 2^{(-10:4)}, C = -5:20
set.seed(2012)
imodel <- train(wlfk~., data = cjb[train set idx, ],</pre>
    method = "svmRadial", preProc = c("center", "scale"),
    tuneGrid = svm grid)
imodel$bestTune
#>
       sigma C
# 121 0.015625 11
```



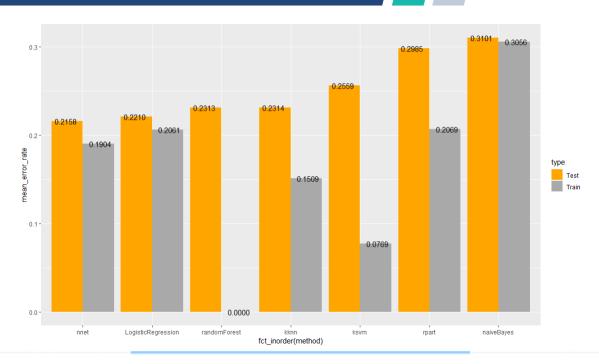
```
sp <- Sys.time() #记录开始时间
cat("\n[Start at:", as.character(sp))
for (i in 1:length(kfolds)) {
  curr fold <- kfolds[[i]] #当前这一折
  train set <- cjb[-curr fold,] #训练集
  test set <- cjb[curr fold,] #测试集
  imodel kfold <- ksvm(wlfk~., data = train set</pre>
      C = imodel$bestTune$C, gamma = imodel$bestTune$sigma)
  predicted train <- predict(imodel kfold,</pre>
           train set, type = "response")
```

```
imetrics("ksvm", "Train",
           predicted train, train set$wlfk)
  predicted test <- predict(imodel kfold,</pre>
           test set, type = "response")
  imetrics("ksvm", "Test",
           predicted test, test set$wlfk)
ep <- Sys.time()
cat("\tFinised at:", as.character(ep), "]\n")
cat("[Time Ellapsed: \t",
    difftime(ep, sp, units = "secs"), " seconds]\n")
```

```
#> method type accuracy error rate
        ksvm Train 0.9181034 0.08189655
#> 1
#> 2
       ksym Test 0.7307692 0.26923077
       ksvm Train 0.9252874 0.07471264
#> 3
#> 4
       ksvm Test 0.7692308 0.23076923
#> 5
        ksym Train 0.9181034 0.08189655
#> 6
       ksvm Test 0.8205128 0.17948718
       ksvm Train 0.9167862 0.08321377
#> 19
#> 20
       ksym Test 0.7662338 0.23376623
```

```
global performance %>%
  group by (method, type) %>%
  summarise(mean error rate = mean(error rate)) %>%
  arrange(type, mean error rate) %>%
  ggplot(aes(x = fct inorder(method), y = mean error rate,
             fill = type)) +
  geom bar(stat= "identity", position = "dodge") +
  geom text(aes(label = format(mean error rate, digits = 3)),
            position = position dodge(width = 1))+
  scale fill manual(values=c("orange","darkgrey")) +
  theme(axis.text.x = element text(angle = 60, hjust = 1))
```

## 模型综合评估



# 謝謝聆听 Thank you

## 教师个人联系方式

## 艾新波

手机: 13641159546

QQ: 23127789

微信: 13641159546

E-mail: 13641159546@126.com

axb@bupt.edu.cn

地址:北京邮电大学科研楼917室

课程 网址: https://github.com/byaxb/RDataAnalytics



