

L6: SVM II

Shan Wang

Lingnan College, Sun Yat-sen University

2020 Data Mining and Machine Learning LN3119

<https://wangshan731.github.io/DM-ML/>

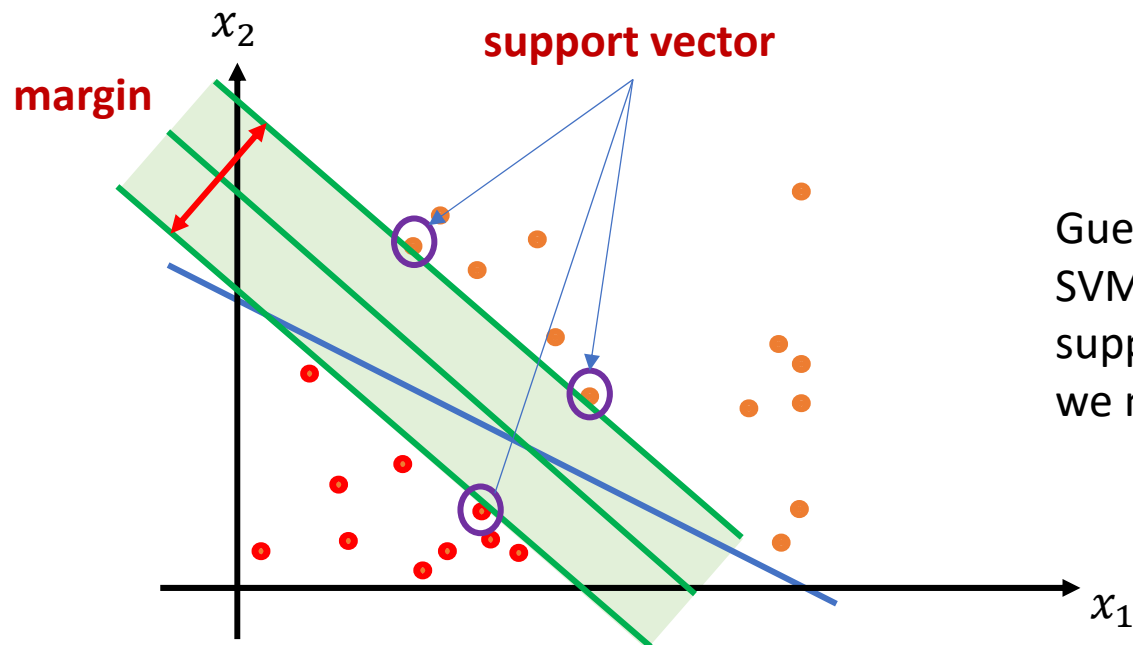


Last and this lecture

- ✓ Linear SVM
 - ✓ Model
 - ✓ Strategy
 - ✓ Algorithm
- ✓ Regularization
 - Kernels
 - Application: diabetes care revisit

Linear SVM

$$y = f_{\theta}(\mathbf{x}) = \begin{cases} +1, & \text{if } \boldsymbol{\theta}'\mathbf{x} + \theta_0 \geq 0 \\ -1, & \text{if } \boldsymbol{\theta}'\mathbf{x} + \theta_0 < 0 \end{cases}$$



Guess: in linear SVM, how many support vectors we need?

Prim

$$\begin{aligned} \min_{\boldsymbol{\theta}, \theta_0} \quad & \frac{1}{2} \|\boldsymbol{\theta}\|^2 \\ \text{s.t.} \quad & y_i(\boldsymbol{\theta}'\mathbf{x}_i + \theta_0) \geq 1, i = 1, \dots, N \end{aligned}$$

Dual

$$\begin{aligned} \max_{\alpha \geq 0} \quad & \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_j' \mathbf{x}_i \\ \text{s.t.} \quad & \sum_{i=1}^N \alpha_i y_i = 0 \end{aligned}$$

Regularization: soft margin

$$\begin{aligned} \min_{\boldsymbol{\theta}, \theta_0} \quad & \frac{1}{2} \|\boldsymbol{\theta}\|^2 + C \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & y_i(\boldsymbol{\theta}'\mathbf{x}_i + \theta_0) \geq 1 - \xi_i, i = 1, \dots, N \\ & \xi_i \geq 0, \quad i = 1, \dots, N \end{aligned}$$

Quadratic programming

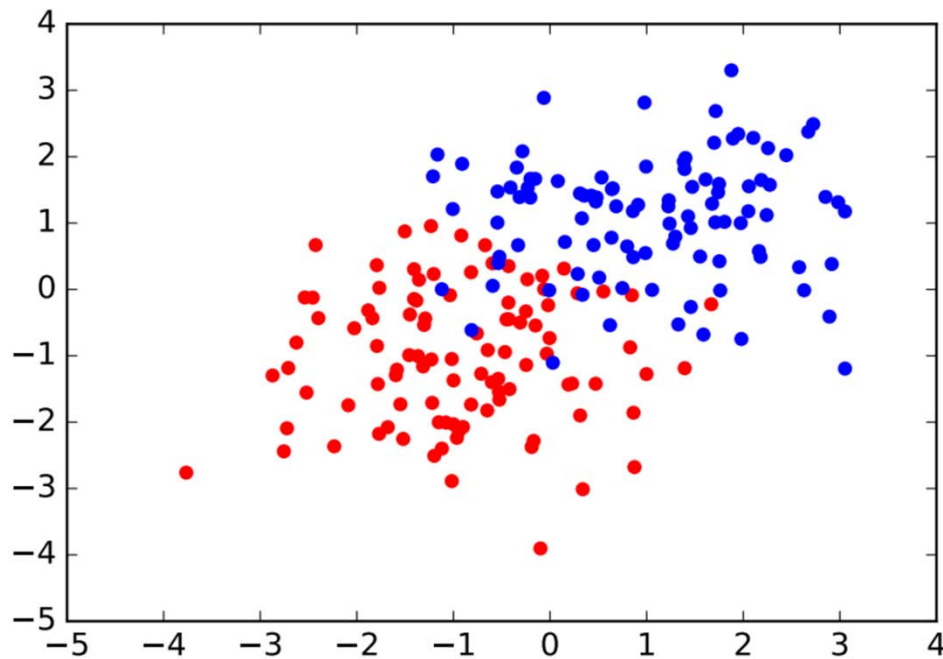
$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_j' \mathbf{x}_i \\ \text{s.t.} \quad & \sum_{i=1}^N \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C, \quad i = 1, \dots, N \end{aligned}$$

SMO algorithm

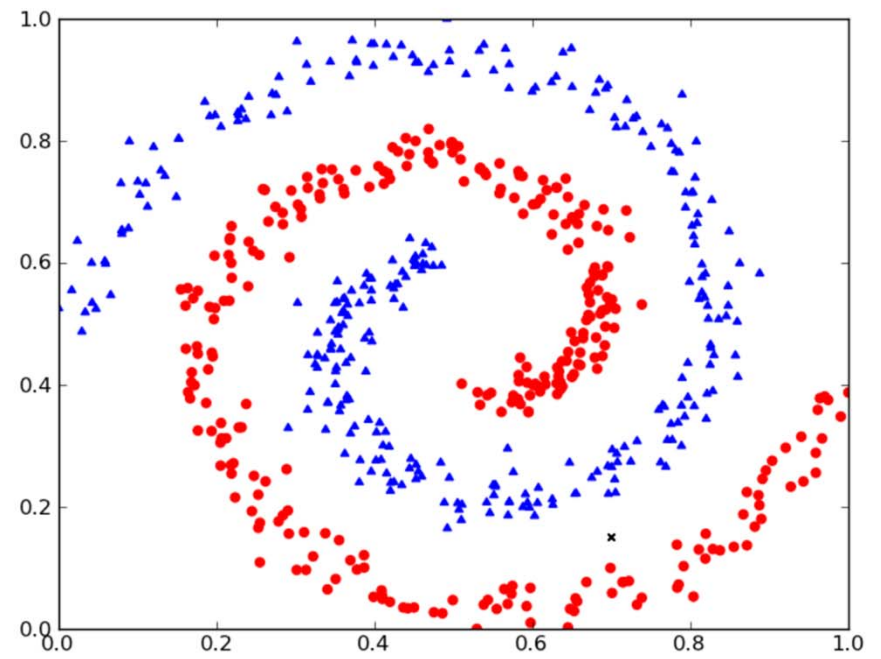
Kernels

Non-separable Data

- When soft margin does not work



May be solved by soft margin

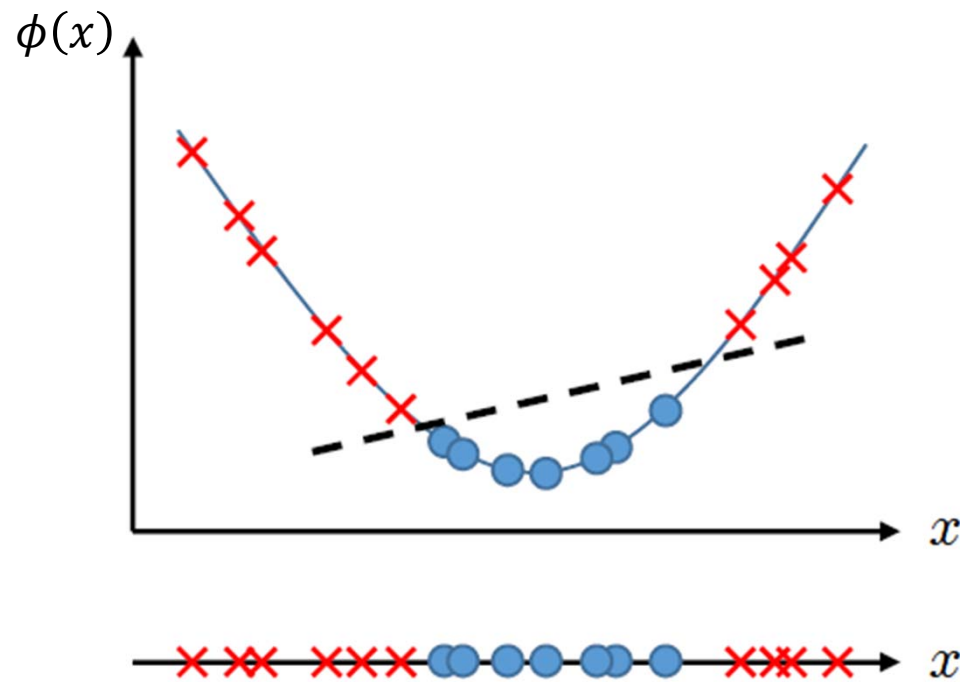


Cannot be solved by soft margin

Non-separable Data (cont.)

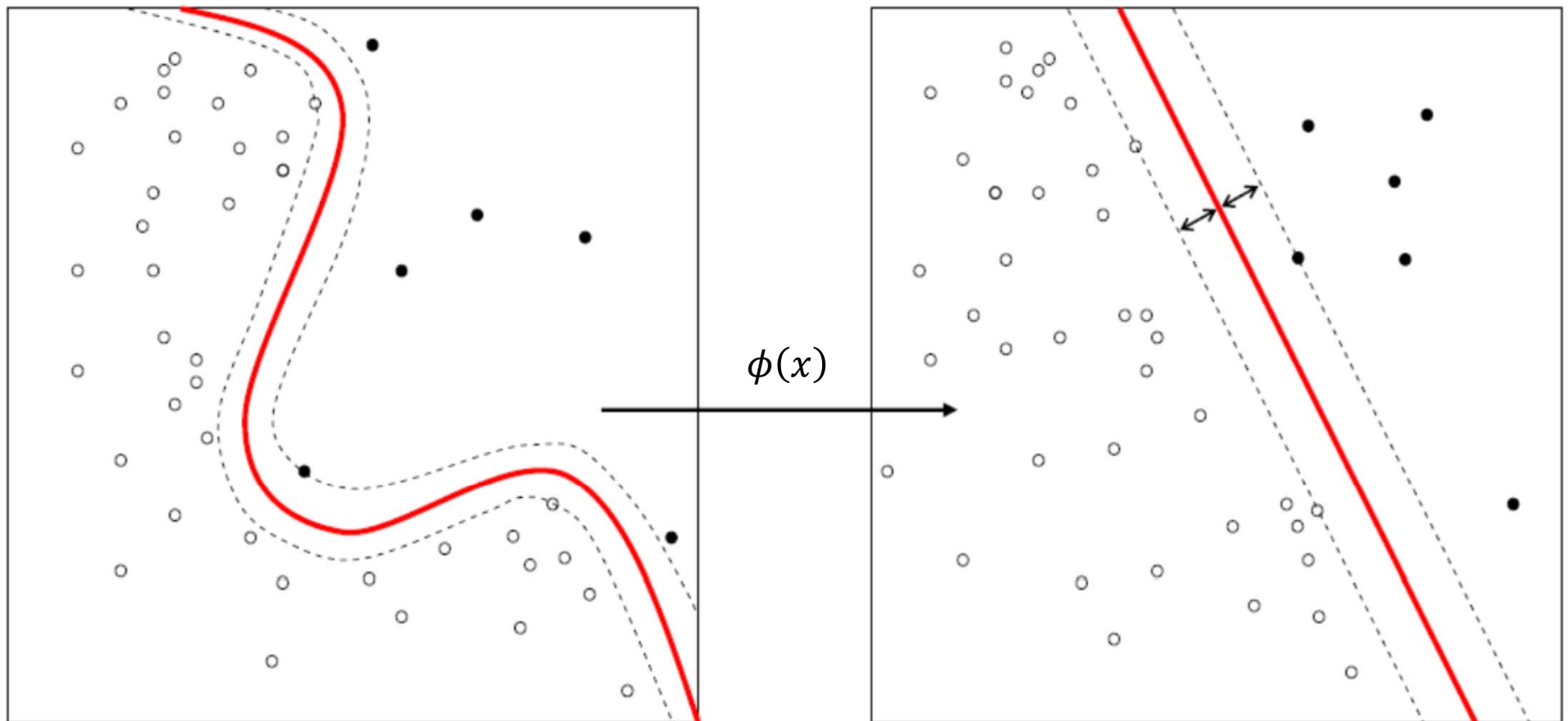
- Mapping feature vectors to a different space
- For example

$$\phi(x) = x^2$$



Non-separable Data (cont.)

- More generally,



From linear SVM to kernel SVM

- Linear SVM:

$$\sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \boxed{\mathbf{x}_i' \mathbf{x}_j}$$

- the inner product of $\mathbf{x}_i, \mathbf{x}_j$: $\mathbf{x}_i' \mathbf{x}_j$

- After feature mapping

$$\sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \boxed{\phi(\mathbf{x}_i)' \phi(\mathbf{x}_j)}$$

- the inner product of $\phi(\mathbf{x}_i), \phi(\mathbf{x}_j)$: $\phi(\mathbf{x}_i)' \phi(\mathbf{x}_j)$

- Kernel function: $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)' \phi(\mathbf{x}_j)$

Kernel trick

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)' \phi(\mathbf{x}_j)$$

- With the example feature mapping function

$$\phi(x) = \begin{bmatrix} x \\ x^2 \\ x^3 \end{bmatrix}$$

- The corresponding kernel is

$$K(x_i, x_j) = \phi(x_i)' \phi(x_j) = x_i x_j + x_i^2 x_j^2 + x_i^3 x_j^3$$

[Trick]

- For many cases, without defining $\phi(\cdot)$, we can directly define $K(\mathbf{x}_i, \mathbf{x}_j)$ for each pair of $\mathbf{x}_i, \mathbf{x}_j$
- For prediction, only need $K(\mathbf{x}_i, \mathbf{x})$ on support vector \mathbf{x}_i

Kernel matrix

- For a finite set of instances $\{x_1, x_2, \dots, x_N\}$
- The kernel matrix K is defined as $[K_{ij}]_{i,j=1,2,\dots,N}$ where $K_{ij} = K(x_i, x_j)$
- If $K(\cdot, \cdot)$ is a valid kernel (that is, is defined by some feature mapping ϕ), then the corresponding kernel matrix $K = [K_{ij}]_{ij} \in \mathbb{R}^{N \times N}$ is **symmetric positive semi-definite** matrix
 - Symmetric: $K_{ij} = K_{ji}$
 - Positive semi-definite: for $\forall z \in \mathbb{R}^N$, we have $z'Kz \geq 0$

How to prove it?

Example valid kernels

- Gaussian kernel
 - $K(\mathbf{x}, \mathbf{z}) = \exp\left(-\frac{\|\mathbf{x}-\mathbf{z}\|^2}{2\sigma^2}\right)$
 - Radial basis function (RBF) kernel
- Simple polynomial kernel
 - $K(\mathbf{x}, \mathbf{z}) = (\mathbf{x}'\mathbf{z})^d$
- Cosine similarity kernel
 - $K(\mathbf{x}, \mathbf{z}) = \frac{\mathbf{x}'\mathbf{z}}{\|\mathbf{x}\|\|\mathbf{z}\|}$
- Sigmoid kernel
 - $K(\mathbf{x}, \mathbf{z}) = \tanh(\alpha \mathbf{x}'\mathbf{z} + c)$
 - $\tanh(b) = \frac{1 - \exp(-2b)}{1 + \exp(-2b)}$

Pros and cons of SVM

- Advantages:

- The solution is globally optimal
 - Based on convex optimization
- Can be applied to both linear/non-linear classification problems
- Can be applied to high-dimensional data
 - The complexity of the data set mainly depends on the support vectors
- Complete theoretical guarantee
 - Compared with deep learning

- Disadvantages:

- The number of parameters α is number of samples, thus hard to apply to large-scale problems
 - SMO can ease the problem a bit
- Mainly applies to binary classification problems
 - For multi-classification problems, can solve several binary classification problems, but might face the problem of imbalanced data

Application

Application: Diabetes Care Revisit

Claims Data

- Predict poor quality care or not

Medical Claims

Diagnosis, Procedures,
Doctor/Hospital, Cost

Pharmacy Claims

Drug, Quantity, Doctor,
Medication Cost

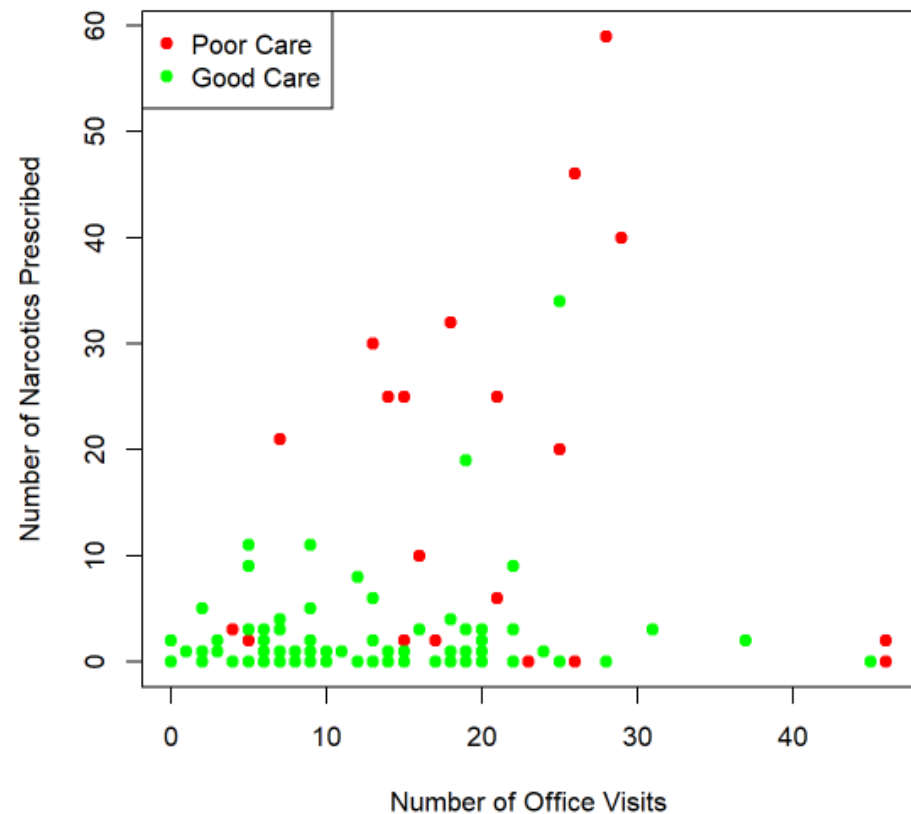
- Electronically available
- Standardized
- Not 100% accurate
- Under-reporting is common
- Claims for hospital visits can be vague

Building a model

- We use a **SVM**
 - Predicts an outcome variable, or *dependent/response variable*
 - Using a set of *independent/explanatory variables*
- Dependent variable: Poor care or not
 - is equal to 1 if the patient had poor care, and equal to -1 if the patient had good care
- Independent variables:
 - Number of Office Visits (OfficeVisits)
 - Number of Narcotics Prescribed (Narcotics)
 - Etc.

Model for Healthcare Quality

- Plot of the independent variables
 - Number of Office Visits (OfficeVisits)
 - Number of Narcotics Prescribed (Narcotics)
- Red are poor care
- Green are good care
- Are these variables predictive of good care or poor care?



Application: Diabetes Care Revisit

Let's get our hands dirty!

Data Analysis

Let's do some basic data analysis using our WHO data.

Hide

```
WHO$under15
```

```
[1] 47.42 23.33 27.42 15.20 47.58 25.96 24.42 20.34 18.95 14.51 22.25 21.62 20.16 30.57 18.99 15.10 16.88 34.4  
0 42.95 28.53  
[21] 35.23 16.35 33.75 24.56 25.75 13.53 45.66 44.20 31.23 43.88 16.37 30.17 40.87 48.52 21.38 17.95 28.03 42.1  
7 42.37 30.61  
[41] 23.94 41.48 14.08 16.58 17.16 14.56 21.98 45.11 17.66 33.73 25.96 30.53 30.29 31.25 30.63 38.95 43.10 15.6  
9 43.29 28.88  
[61] 16.42 18.26 38.49 45.90 17.62 13.17 38.59 14.60 26.96 40.80 42.46 41.55 36.77 35.35 35.72 14.62 20.71 29.4  
3 29.27 23.68  
[81] 40.51 23.54 27.53 14.84 37.78 13.13 34.13 25.46 42.37 30.10 24.90 30.21 35.61 14.57 21.64 36.75 43.05 29.4  
5 15.13 17.46  
[101] 42.73 45.44 26.65 29.83 47.14 14.98 30.10 40.22 20.17 29.82 35.81 18.26 27.85 19.81 27.85 45.38 25.28 36.5  
9 30.10 35.58  
[121] 17.21 20.26 33.37 49.99 44.23 30.61 18.64 24.19 34.31 30.10 28.65 38.37 32.78 29.18 34.53 14.91 14.92 13.2  
8 15.25 16.52  
[141] 15.05 15.45 43.56 25.96 24.31 25.70 37.88 14.04 41.60 29.69 43.54 16.45 21.95 41.74 16.48 15.00 14.16 40.3  
7 47.35 29.53  
[161] 42.28 15.20 25.15 41.48 27.83 38.05 16.71 14.79 35.35 35.75 18.47 16.89 46.33 41.89 37.33 20.73 23.22 26.8  
0 28.65 30.61  
[181] 48.54 14.18 14.41 17.54 44.85 19.63 22.05 28.90 37.37 28.84 22.87 40.72 46.73 40.34
```

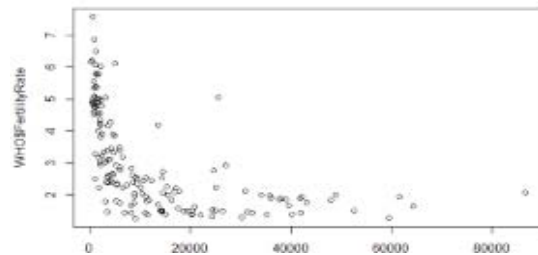
```
WHO$country[which.min(WHO$under15)]
```

```
[1] Japan  
294 Levels: Afghanistan Albania Algeria Andorra Angola Antigua and Barbuda Argentina Armenia Australia Austria  
... Zimbabwe
```

Let's create some plots for exploratory data analysis (EDA). First, let's create a basic scatterplot of GNI versus FertilityRate.

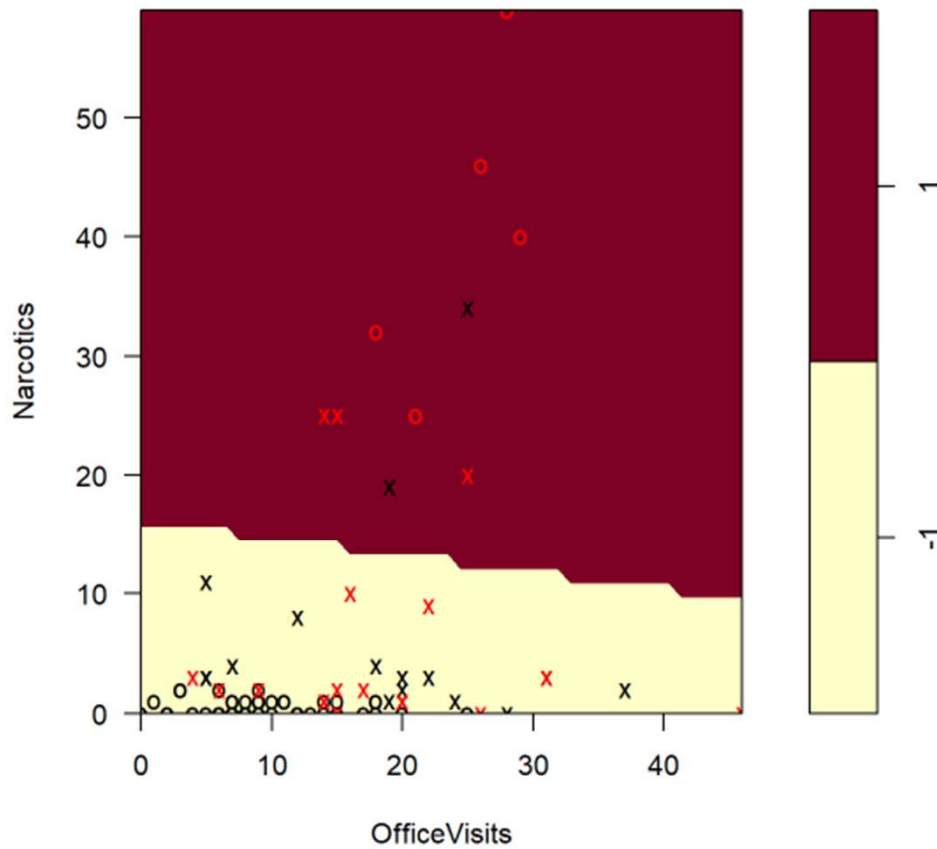
Hide

```
plot(WHO$GNI, WHO$fertilityRate)
```



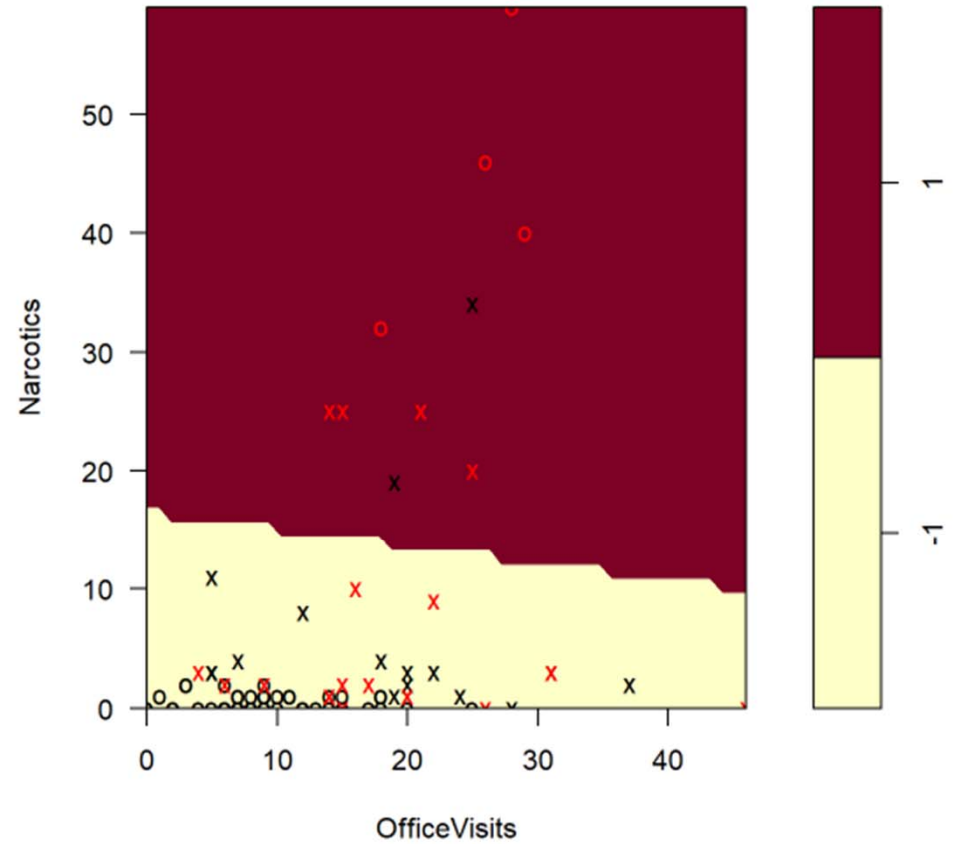
Linear SVM

SVM classification plot



Large C

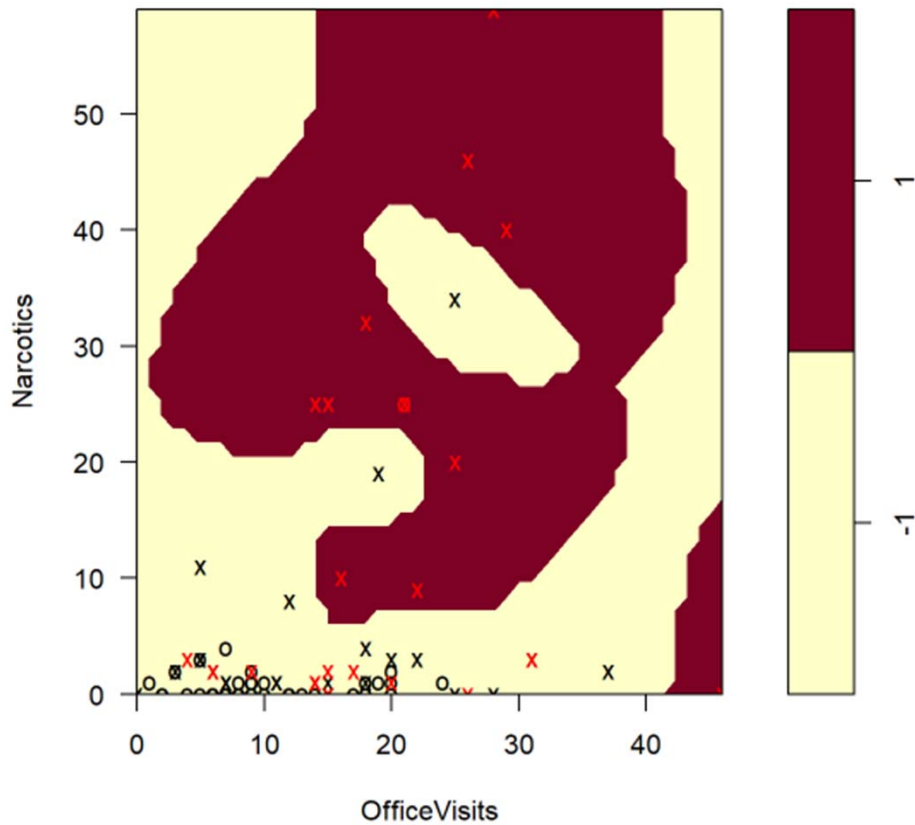
SVM classification plot



Small C

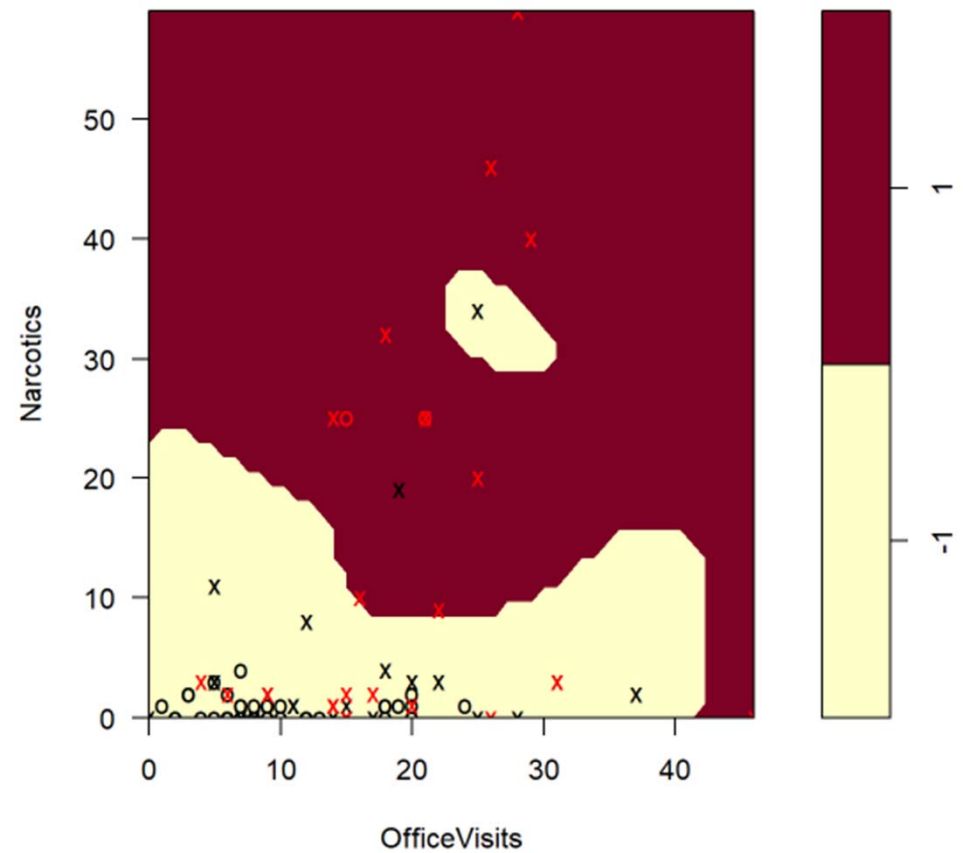
SVM with RBF kernel

SVM classification plot



Large C

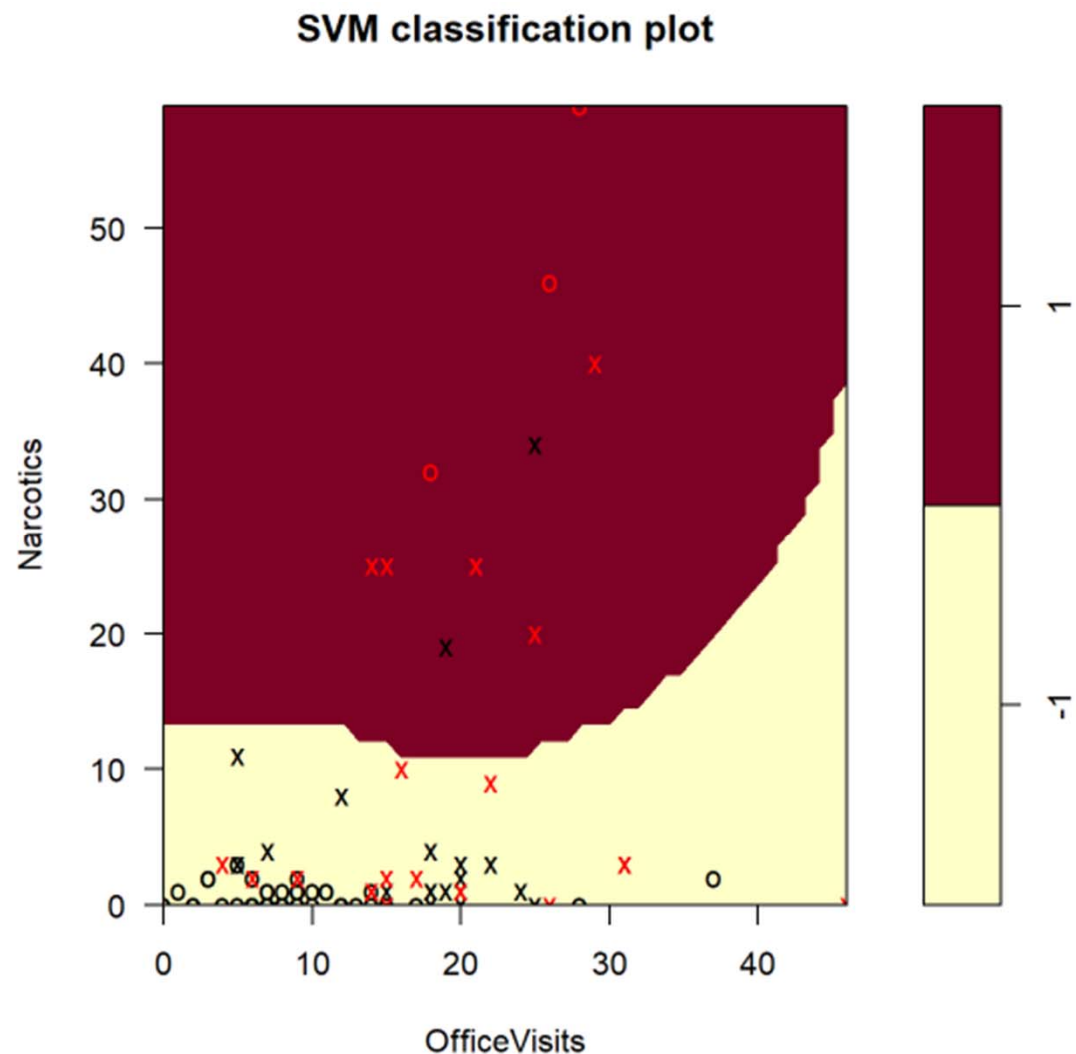
SVM classification plot



Small C

K-fold cross validation and testing

- Best model after tuning
 - Sigmoid kernel
- In-sample accuracy
 - 82.16%
- Out-of-sample
 - 79.49%



Lecture 6 wrap-up

- ✓ Linear SVM
 - ✓ Model
 - ✓ Strategy
 - ✓ Algorithm
- ✓ Regularization
- ✓ Kernels
- ✓ Application: diabetes care revisit

Next lecture

- Supervised learning

- Linear regression
- Logistic regression
- SVM and kernel
- Tree models

- Deep learning

- Neural networks
- Convolutional NN
- Recurrent NN

- Unsupervised learning

- Clustering
- PCA
- EM

- Reinforcement learning

- MDP
- ADP
- Deep Q-Network

2020 Data Mining and Machine Learning LN3119

<https://wangshan731.github.io/DM-ML/>



Questions?

Shan Wang (王杉)

<https://wangshan731.github.io/>