

Kernel methods

Victor Kitov

`v.v.kitov@yandex.ru`

Kernel trick

Perform feature transformation: $x \rightarrow \phi(x)$. Scalar product becomes $\langle x, x' \rangle \rightarrow \langle \phi(x), \phi(x') \rangle = K(x, x')$

Kernel trick

Define not the feature representation x but only scalar product function $K(x, x')$

- Comments:
 - required that the solution depends only on scalar products. Kernels can be constructed from other kernels, for example from:
 - 1 scalar product $\langle x, x' \rangle$
 - 2 constant $K(x, x') \equiv 1$
 - 3 $x^T A x$ for any $A \succcurlyeq 0^1$
 - feature representation $\phi(x)$ not needed
 - $\langle x, x' \rangle$ has complexity $O(D)$. Complexity of $K(x, x')$ may be $O(1)$.

Kernelizable algorithms

- ridge regression:
- K-NN
- K-means
- PCA
- SVM
- many more...

Kernel trick use cases

- high-dimensional data
 - polynomial of order up to M
 - Gaussian kernel $K(x, x') = e^{-\frac{1}{2\sigma^2} \|x-x'\|^2}$ corresponds to infinite-dimensional feature space.
- hard to vectorize data
 - strings, sets, images, texts, graphs, 3D-structures, sequences, etc.
- natural scalar product exist
 - strings: number of co-occurring substrings
 - sets: size of intersection of sets
 - example: for sets S_1 and S_2 : $K(S_1, S_2) = 2^{|S_1 \cap S_2|}$ is a possible kernel.
 - etc.
- scalar product can be computed efficiently

General motivation for kernel trick

- perform generalization of linear methods to non-linear case
 - as efficient as linear methods
 - local minimum is global minimum
 - no local optima=>less overfitting
- non-vectorial objects
 - hard to obtain vector representation

Kernel definition

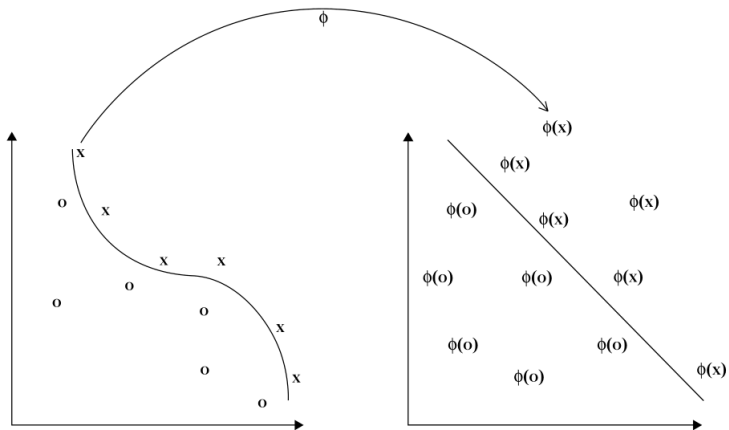
- x is replaced with $\phi(x)$
 - Example: $[x] \rightarrow [x, x^2, x^3]$

Kernel

Function $K(x, x') : X \times X \rightarrow \mathbb{R}$ is a kernel function if it may be represented as $K(x, x') = \langle \phi(x), \phi(x') \rangle$ for some mapping $\phi : X \rightarrow H$, with scalar product defined on H .

- $\langle x, x' \rangle$ is replaced by $\langle \phi(x), \phi(x') \rangle = K(x, x')$

Illustration



Specific types of kernels

- $K(x, x') = K(x - x')$ - stationary kernels (invariant to translations)
- $K(x, x') = K(\|x - x'\|)$ - radial basis functions

Polynomial kernel²

- Example 1: let $D = 2$.

$$\begin{aligned}K(x, z) &= (x^T z)^2 = (x_1 z_1 + x_2 z_2)^2 = \\&= x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 z_1 x_2 z_2 \\&= \phi^T(x) \phi(z)\end{aligned}$$

for $\phi(x) = (x_1^2, x_2^2, \sqrt{2}x_1 x_2)$

²What kind of feature transformation will correspond to $K(x, z) = (x^T z)^M$ for arbitrary M and D ?

Polynomial kernel³

- Example 2: let $D = 2$.

$$\begin{aligned}
 K(x, z) &= (1 + x^T z)^2 = (1 + x_1 z_1 + x_2 z_2)^2 = \\
 &= 1 + x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 z_1 + 2x_2 z_2 + 2x_1 z_1 x_2 z_2 \\
 &= \phi^T(x) \phi(z)
 \end{aligned}$$

for $\phi(x) = (1, x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2)$

³What kind of feature transformation will correspond to $K(x, z) = (1 + x^T z)^M$ kernels for arbitrary M and D ?

Kernel properties

Theorem (Mercer): Function $K(x, x')$ is a kernel is and only if

- it is symmetric: $K(x, x') = K(x', x)$
- it is non-negative definite:
 - definition 1: for every function $g : X \rightarrow \mathbb{R}$

$$\int_X \int_X K(x, x') g(x) g(x') dx dx' \geq 0$$

- definition 2 (equivalent): for every finite set x_1, x_2, \dots, x_M
Gramm matrix $\{K(x_i, x_j)\}_{i,j=1}^M \succeq 0$ (p.s.d.)

Kernel construction

- Kernel learning - separate field of study.
- Hard to prove non-negative definiteness of kernel in general.
- Kernels can be constructed from other kernels, for example from:
 - 1 scalar product $\langle x, x' \rangle$
 - 2 constant $K(x, x') \equiv 1$
 - 3 $x^T A x$ for any $A \succeq 0$ ⁴

⁴Under what feature transformation will case 1 transform to cases 2 and 3? You may use Cholesky decomposition.

Constructing kernels from other kernels

If $K_1(x, x')$, $K_2(x, x')$ are arbitrary kernels, $c > 0$ is a constant, $q(\cdot)$ is a polynomial with non-negative coefficients, $h(x)$ and $\varphi(x)$ are arbitrary functions $\mathcal{X} \rightarrow \mathbb{R}$ and $\mathcal{X} \rightarrow \mathbb{R}^M$ respectively, then these are valid kernels⁵:

- ① $K(x, x') = cK_1(x, x')$
- ② $K(x, x') = K_1(x, x')K_2(x, x')$
- ③ $K(x, x') = K_1(x, x') + K_2(x, x')$
- ④ $K(x, x') = K_1(\varphi(x), \varphi(x'))$
- ⑤ $K(x, x') = h(x)K_1(x, x')h(x')$
- ⑥ $K(x, x') = e^{K_1(x, x')}$

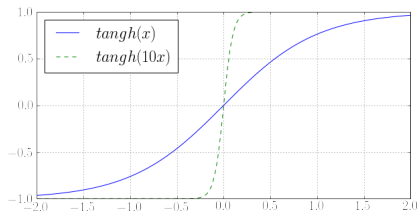
⁵prove some of these statements

Commonly used kernels

Let x and x' be two objects.

Kernel	Mathematical form
linear	$\langle x, x' \rangle$
polynomial	$(\gamma \langle x, x' \rangle + r)^d$
RBF	$\exp(-\gamma \ x - x'\ ^2)$

- Standard transformation is also sigmoid=tangh($\gamma \langle x, y \rangle + r$) but its not a Mercer kernel.



Addition⁶

- Other kernelized algorithms: K-NN, K-means, K-medoids, nearest medoid, PCA, SVM, etc.
- Kernelization of distance:

⁶How can we calculate scalar product between normalized (unit norm) vectors $\phi(x)$ and $\phi(x')$?

Addition⁶

- Other kernelized algorithms: K-NN, K-means, K-medoids, nearest medoid, PCA, SVM, etc.
- Kernelization of distance:

$$\begin{aligned}\rho(\mathbf{x}, \mathbf{x}')^2 &= \langle \phi(\mathbf{x}) - \phi(\mathbf{x}'), \phi(\mathbf{x}) - \phi(\mathbf{x}') \rangle \\ &= \langle \phi(\mathbf{x}), \phi(\mathbf{x}) \rangle + \langle \phi(\mathbf{x}'), \phi(\mathbf{x}') \rangle - 2\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle \\ &= K(\mathbf{x}, \mathbf{x}) + K(\mathbf{x}', \mathbf{x}') - 2K(\mathbf{x}, \mathbf{x}')\end{aligned}$$

⁶How can we calculate scalar product between normalized (unit norm) vectors $\phi(\mathbf{x})$ and $\phi(\mathbf{x}')$?

Table of Contents

- 1 Kernel support vector machines
- 2 Kernel ridge regression

Linear SVM reminder

- Solution for weights:

$$\mathbf{w} = \sum_{i \in \mathcal{SV}} \alpha_i y_i \mathbf{x}_i$$

Discriminant function

$$g(\mathbf{x}) = \sum_{i \in \mathcal{SV}} \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle + w_0$$

$$w_0 = \frac{1}{n_{\tilde{\mathcal{SV}}}} \left(\sum_{j \in \tilde{\mathcal{SV}}} y_j - \sum_{j \in \tilde{\mathcal{SV}}} \sum_{i \in \mathcal{SV}} \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x}_j \rangle \right)$$

where $\mathcal{SV} = \{i : y_i(\mathbf{x}_i^T \mathbf{w} + w_0) \leq 1\}$ are indexes of all support vectors and $\tilde{\mathcal{SV}} = \{i : y_i(\mathbf{x}_i^T \mathbf{w} + w_0) = 1\}$ are boundary support vectors.

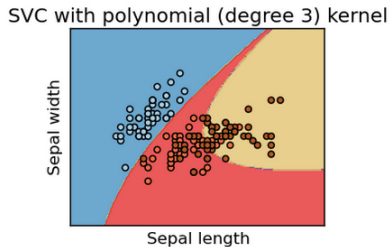
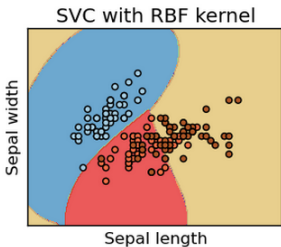
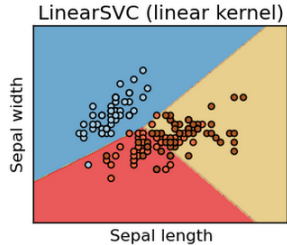
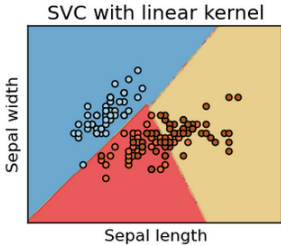
Kernel SVM

Discriminant function

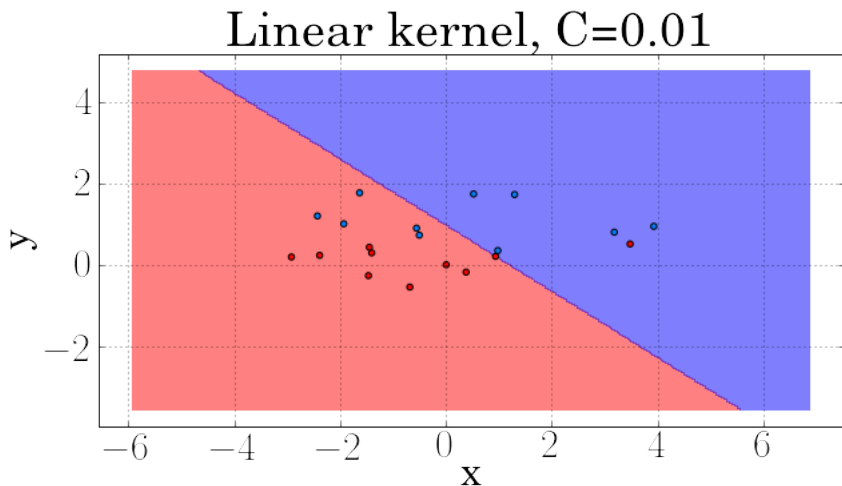
$$g(x) = \sum_{i \in S_V} \alpha_i y_i K(x_i, x) + w_0$$

$$w_0 = \frac{1}{n_{\tilde{S}_V}} \left(\sum_{j \in \tilde{S}_V} y_j - \sum_{j \in \tilde{S}_V} \sum_{i \in S_V} \alpha_i y_i K(x_i, x_j) \right)$$

Kernel results

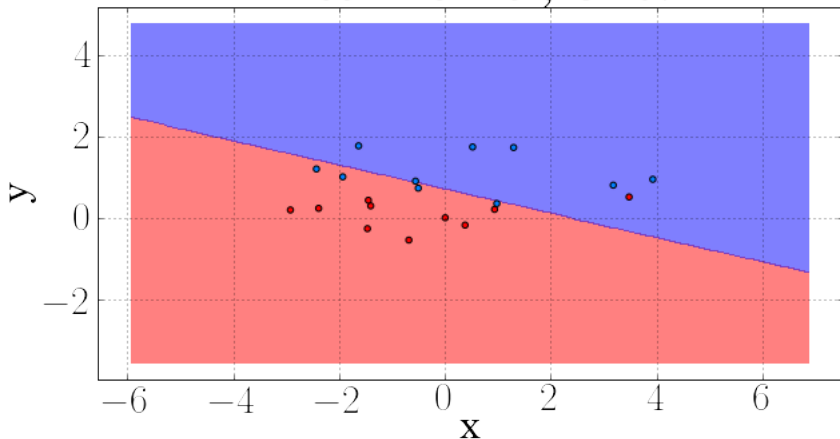


Linear kernel - variable C



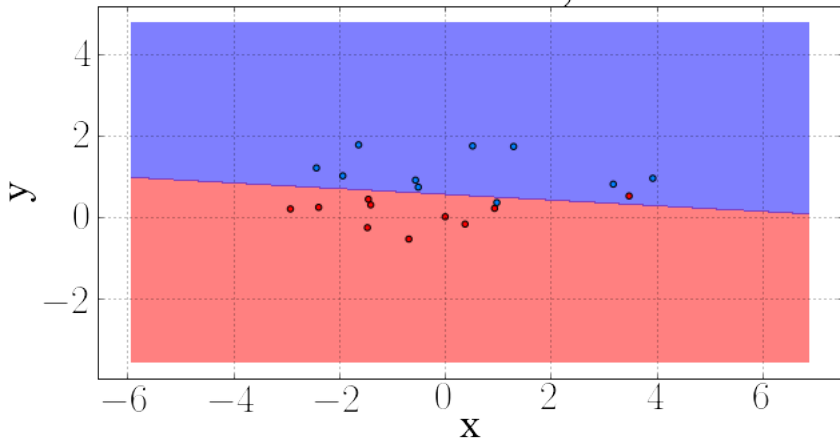
Linear kernel - variable C

Linear kernel, $C=0.1$



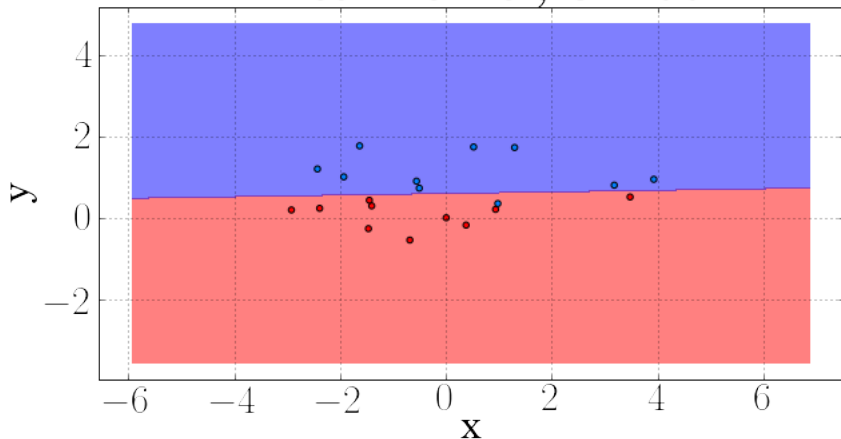
Linear kernel - variable C

Linear kernel, $C=1$

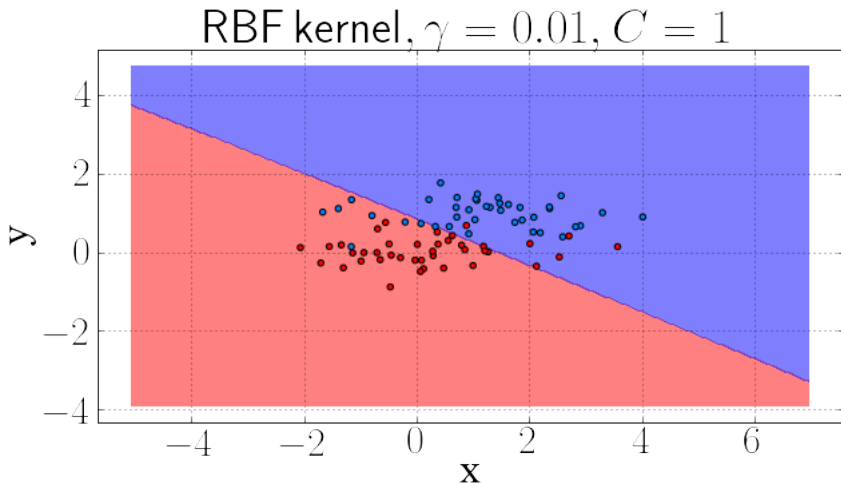


Linear kernel - variable C

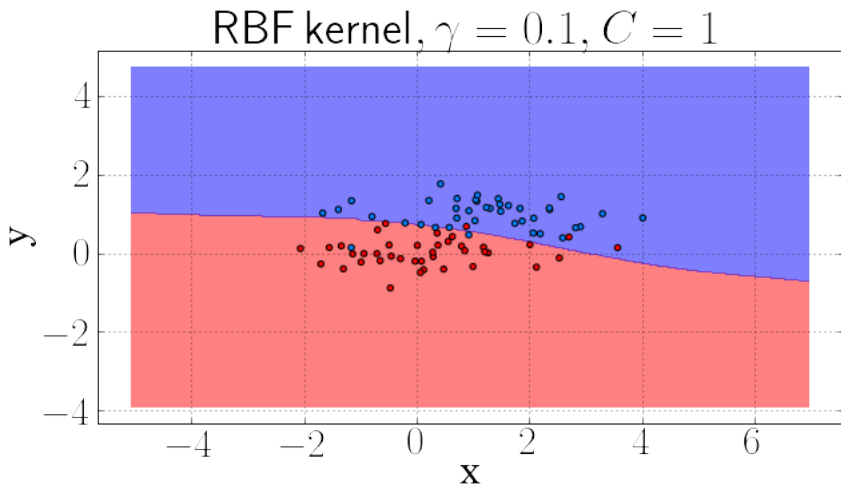
Linear kernel, $C=100$



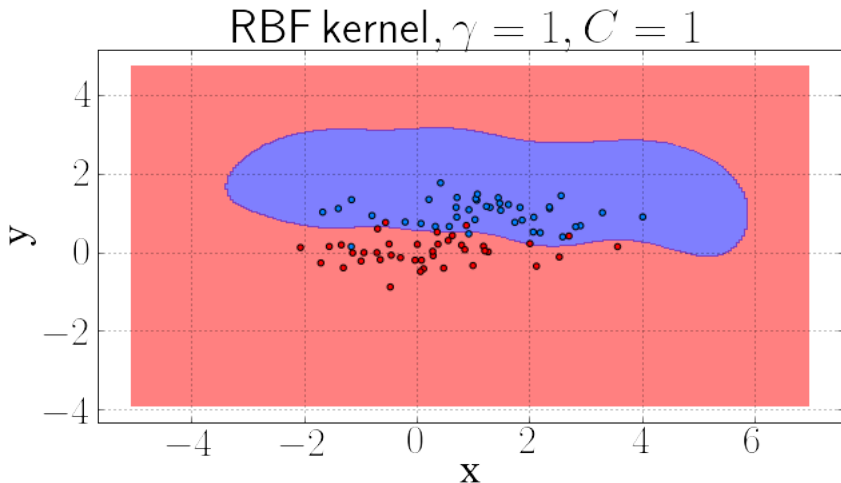
RBF kernel - variable γ



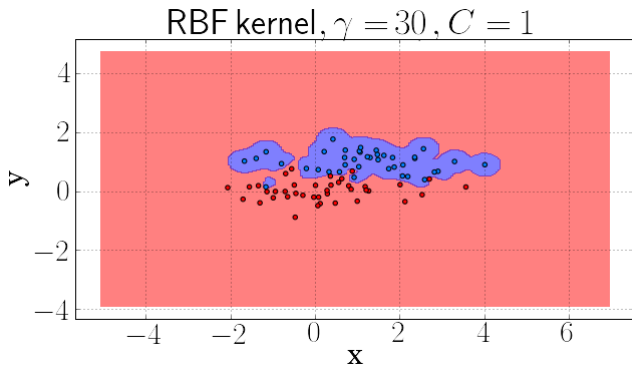
RBF kernel - variable γ



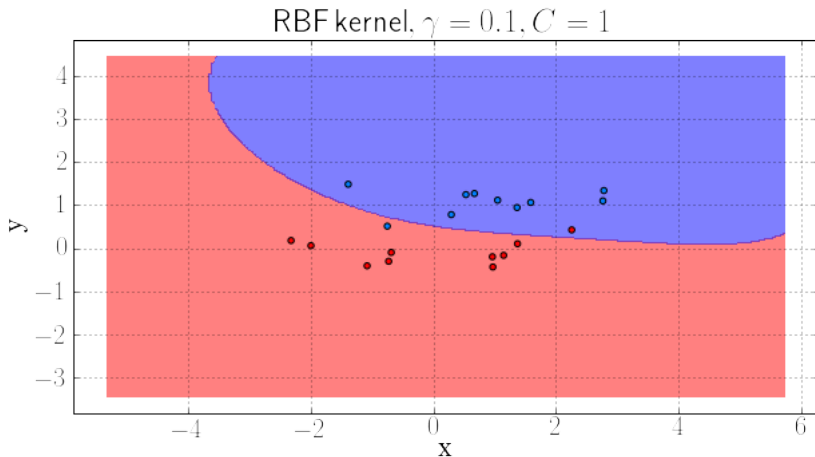
RBF kernel - variable γ



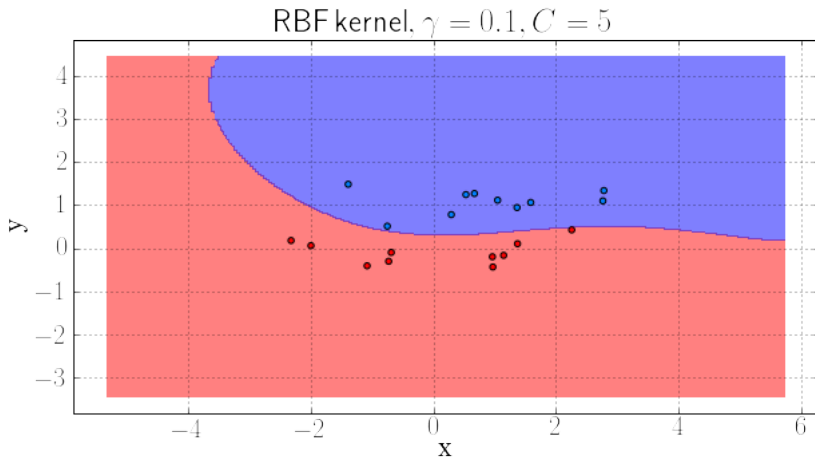
RBF kernel - variable γ



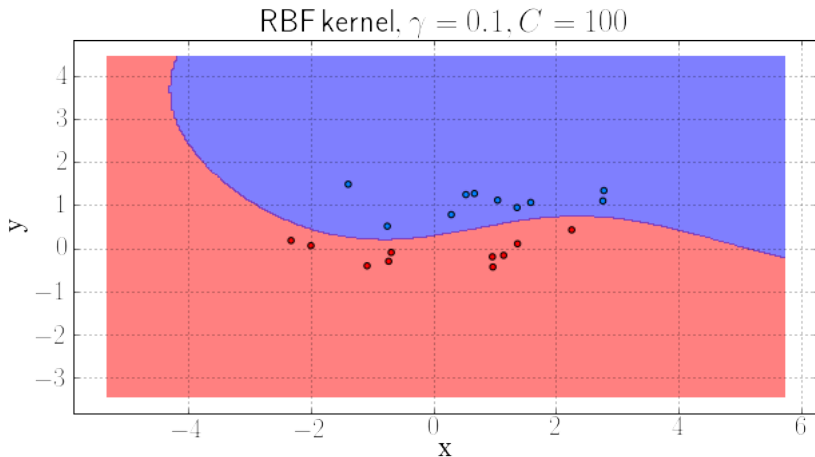
RBF kernel - variable C



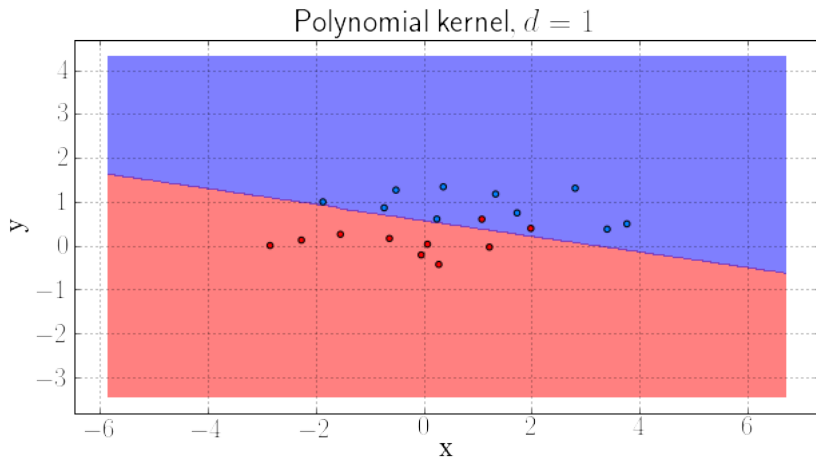
RBF kernel - variable C



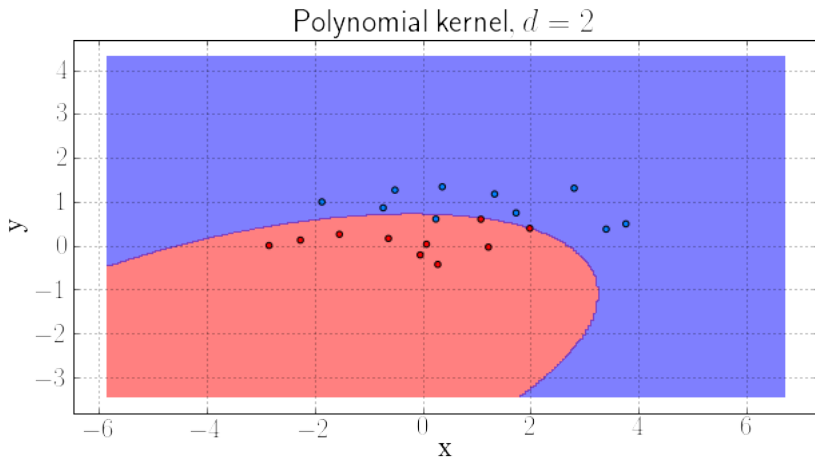
RBF kernel - variable C



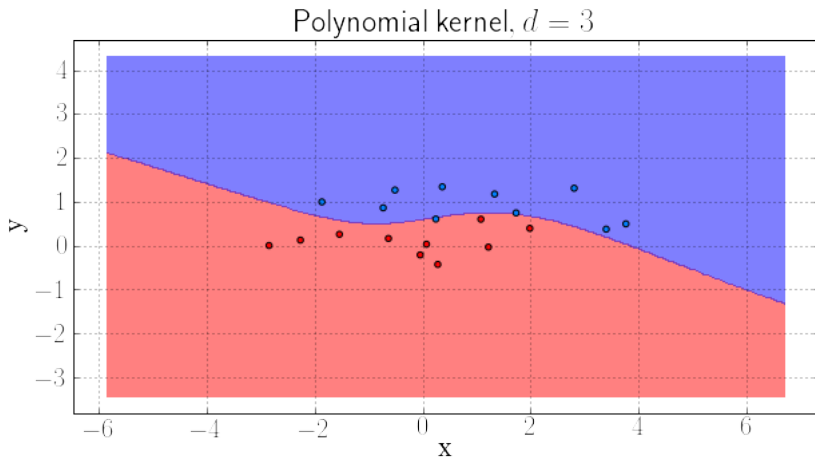
Polynomial kernel - variable d



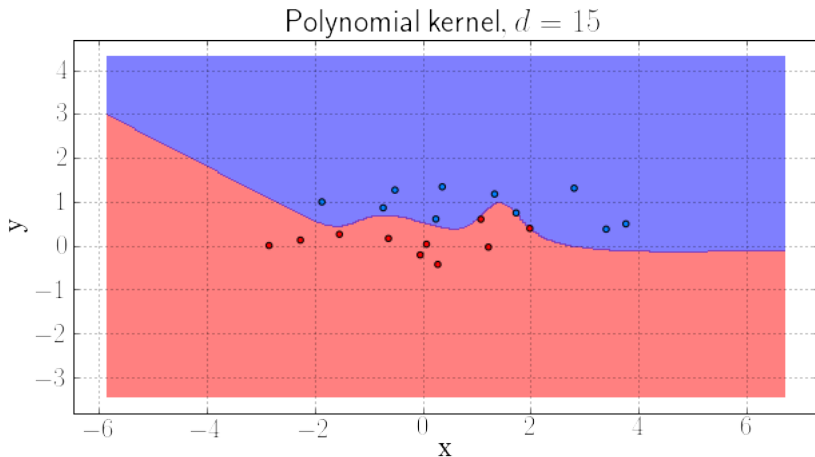
Polynomial kernel - variable d



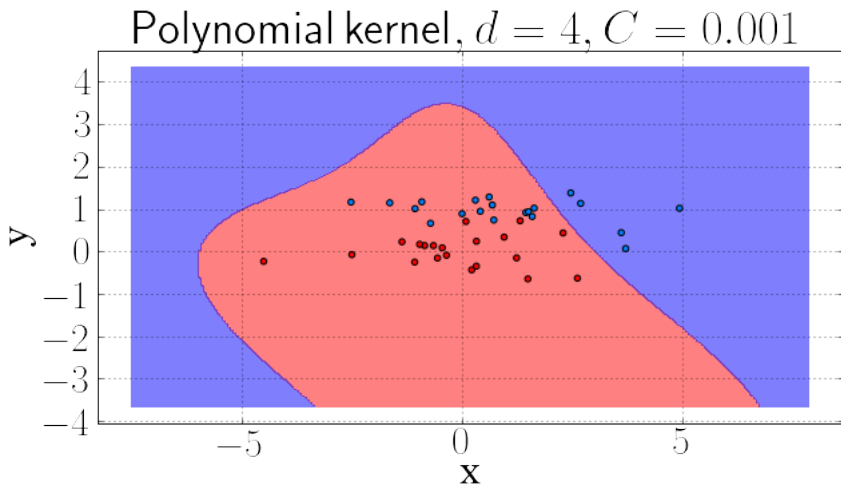
Polynomial kernel - variable d



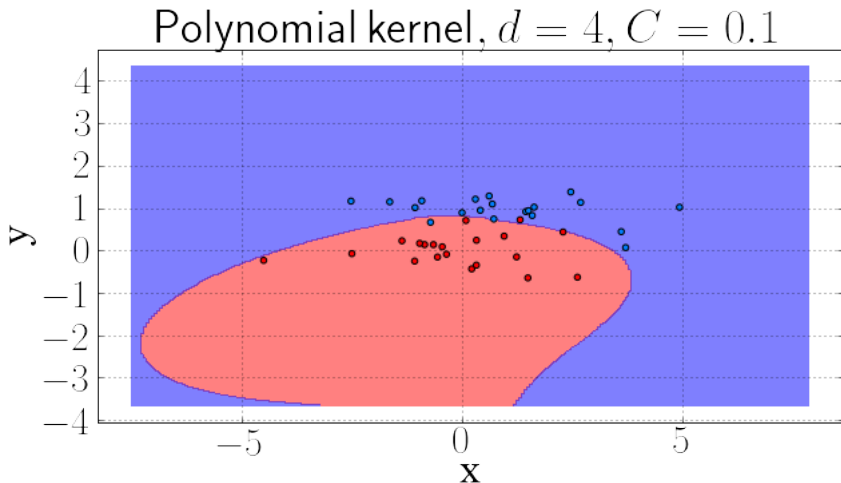
Polynomial kernel - variable d



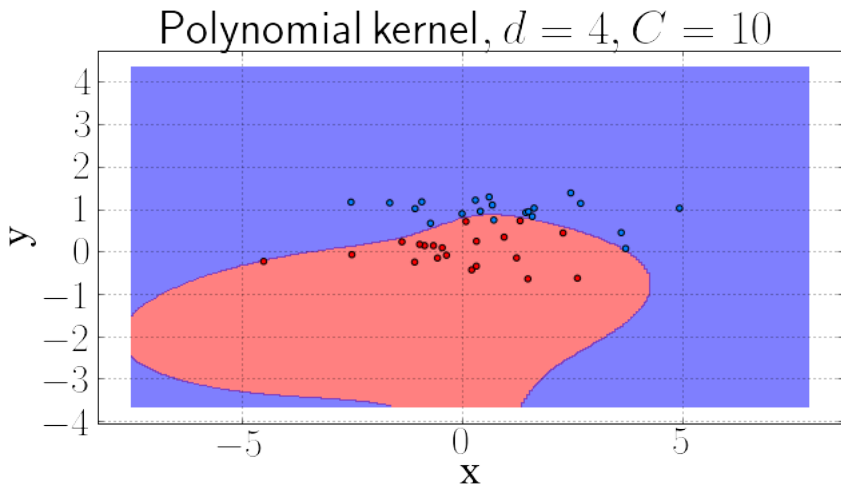
Polynomial kernel - variable C



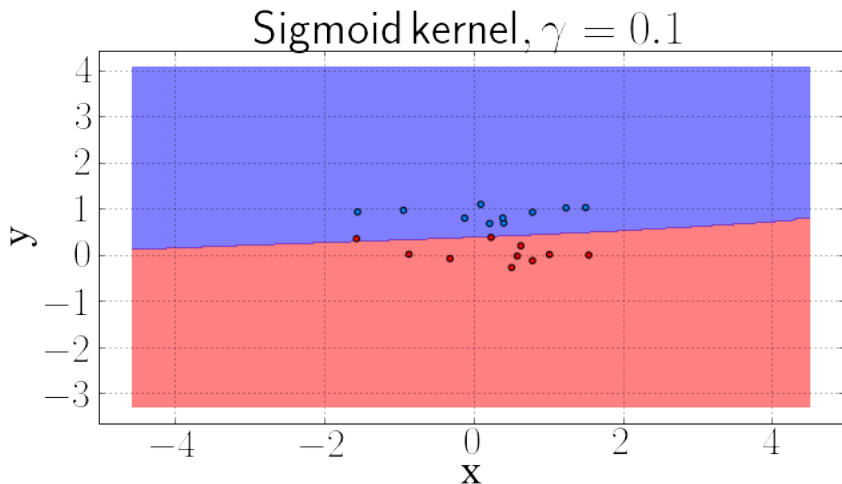
Polynomial kernel - variable C



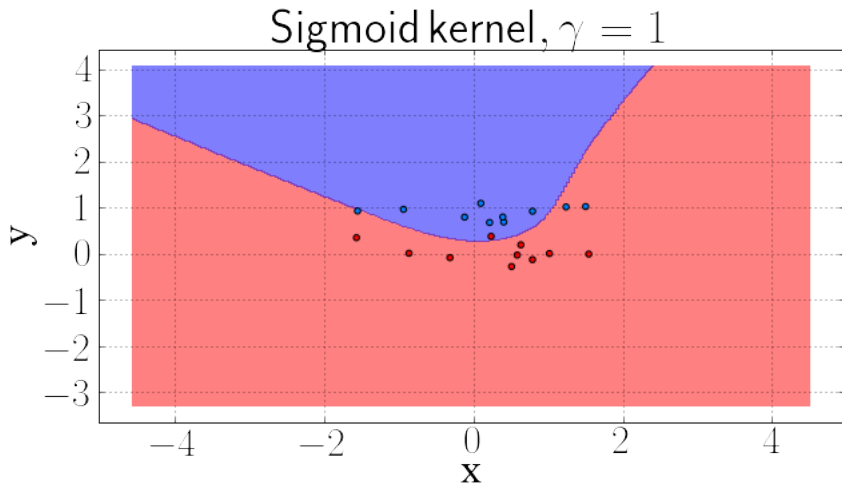
Polynomial kernel - variable C



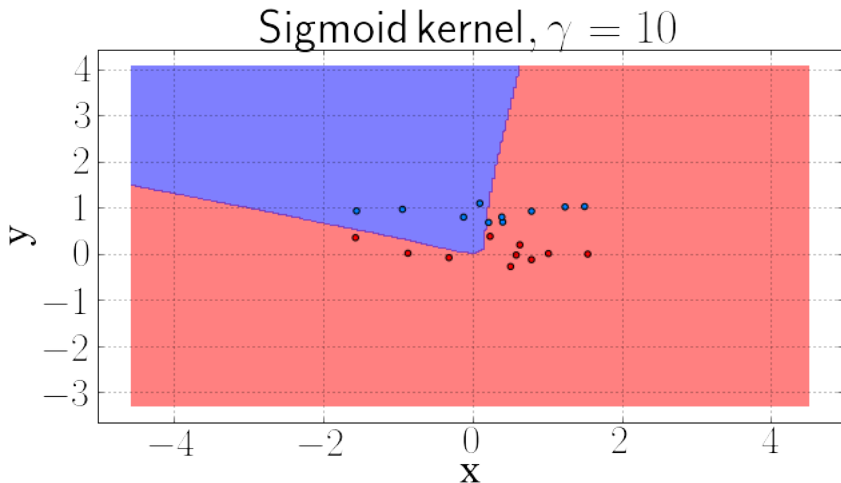
Sigmoid kernel - variable γ



Sigmoid kernel - variable γ



Sigmoid kernel - variable γ



Sigmoid kernel - variable C

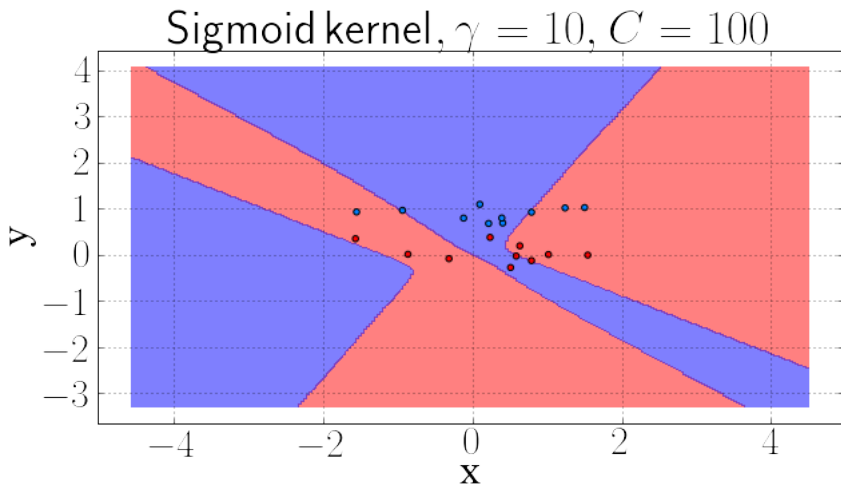


Table of Contents

- 1 Kernel support vector machines
- 2 Kernel ridge regression**

Ridge regression

- Ridge regression criterion:

$$Q(\beta) = \sum_{n=1}^N \left(x_n^T \beta - y_n \right)^2 + \lambda \sum_{d=1}^D \beta_d^2 \rightarrow \min_{\beta}$$

- Stationarity condition:

$$\frac{dQ(\beta)}{d\beta} = 2 \sum_{n=1}^N \left(x_n^T \beta - y_n \right) x_n + 2\lambda\beta = 0$$

- In vector form:

$$X^T (X\beta - Y) + \lambda\beta = 0$$

Ridge regression

- Primal solution:

$$X^T X + \lambda I \beta = X^T Y$$

$$\beta = (X^T X + \lambda I)^{-1} X^T Y$$

- Comment: $X^T X \succcurlyeq 0$ (positive semi-definite) and $X^T X + \lambda I \succ 0$ (positive definite), so ridge regression is always identifiable.
- Cost of estimation:
 - $X^T X + \lambda I$: $ND^2 + D$
 - $X^T Y$: DN
 - $(X^T X + \lambda I)^{-1}$: D^3
 - $(X^T X + \lambda I)^{-1} X^T Y$: D^2
 - Total training cost is $O(ND^2 + D^3) = O(D^2(N + D))$.
- Cost of prediction $\hat{y}(x) = \langle x, \beta \rangle$ is D .

Dual solution

From vector stationarity condition:

$$\mathbf{X}^T (\mathbf{X}\beta - \mathbf{Y}) + \lambda\beta = \mathbf{0}$$

follows the dual solution (a linear combination of training vectors):

$$\beta = \frac{1}{\lambda} \mathbf{X}^T (\mathbf{Y} - \mathbf{X}\beta) = \mathbf{X}^T \alpha \quad (1)$$

where

$$\alpha = \frac{1}{\lambda} (\mathbf{Y} - \mathbf{X}\beta) \quad (2)$$

is called a vector of *dual variables*.

Prediction:

$$\hat{y}(x) = x^T \beta = x^T \mathbf{X}^T \alpha = \sum_{i=1}^N \alpha_i \langle x, x_i \rangle$$

Dual solution

To find α we plug (1) into (2):

$$\begin{aligned}\alpha &= \frac{1}{\lambda}(Y - X\beta) = \frac{1}{\lambda}(Y - XX^T\alpha) \\ (XX^T + \lambda I)\alpha &= Y \\ \alpha &= (XX^T + \lambda I)^{-1}Y\end{aligned}$$

Cost of estimation:

$$XX^T + \lambda I: N^2D + N$$

$$(XX^T + \lambda I)^{-1}: N^3$$

$$(XX^T + \lambda I)^{-1}Y: N^2$$

Total training cost is $O(N^2D + N^3) = O(N^2(D + N))$.

Cost of prediction $\hat{y}(x) = \langle x, \beta \rangle$ is ND .

Dual solution motivation

- Optimal α depends not on exact features but only on scalar products:

$$\alpha = (XX^T + \lambda I)^{-1} Y = (G + \lambda I)^{-1} Y$$

where $G \in \mathbb{R}^{N \times N}$ and $\{G\}_{ij} = \langle x_i, x_j \rangle$ - G is called *Gram matrix*.

- Prediction also depends only on scalar products:

$$\hat{y}(x) = \sum_{i=1}^N \alpha_i \langle x, x_i \rangle = \alpha^T v$$

where $v \in \mathbb{R}^N$ and $v_i = \langle x, x_i \rangle$.