

Support Vector Machines mit Kernen

Jonas Krug, Hendrik Sieck und Tim Schlottmann
TU Hamburg

30. Januar 2018

1 Grundlagen

Support Vector Maschinen sind ein binärer Klassifizierer. Datenmengen werden also in zwei Klassen eingeteilt. Zur Unterteilung der Klassen wird eine Hyperebene benutzt. Wie die Unterteilung stattfindet soll in diesem Abschnitt erklärt werden.

Folgende Definitionen bilden die Grundlage für SVMs:

- Anzahl an Datenpunkten: $m \in \mathbb{R}$
- Input: $\mathbf{x} \in \mathbb{R}^N$
- Output: $y \in \{-1, +1\}$
- Trainingsset: $S \in (\mathbb{R}^N \times \{-1, +1\})^m$

Ziel ist es, eine Hypothese zu finden, die einen Input \mathbf{x} auf eine der beiden Klassen y abbildet:

$$h : \mathbb{R}^N \rightarrow \{+1, -1\}$$
$$\mathbf{x} \mapsto y$$

1.1 Hyperebene

Betrachten wir ein Trainingsset S mit eingezeichneter Hyperebene H wie in Abbildung 1. Es fällt auf, dass ein Spalt (engl. margin) entsteht, dessen Grenzen sich durch zwei weitere Hilfsebenen H_+ und H_- darstellen lassen. H_+ und H_- sind hierbei parallel zu H . In Abbildung 1 ist H_+ die Begrenzung zur $+1$ -Klasse und H_- die Begrenzung zur -1 -Klasse. Die Vektoren, durch die die Hilfsebenen verlaufen, heißen Supportvektoren und sind in der Abbildung schwarz umrandet.

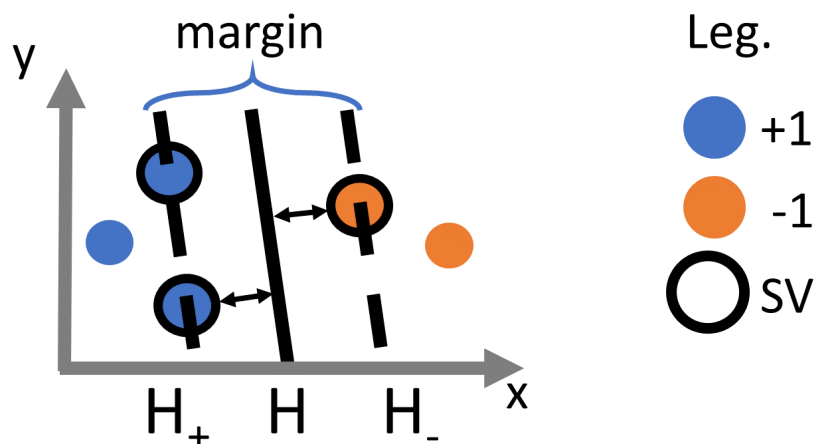


Abbildung 1: Hyperebene H zur Trennung der beiden Klassen $+1$ und -1

Nun gilt es die Hypothese zu finden. Hierfür betrachten wir zuerst die Hyperbenengleichung $H = \mathbf{w}^T \mathbf{x} + b = 0$. Hierbei werden die Parameter \mathbf{w} und b so gewählt, dass für die Supportvektoren gilt: $y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1$. So lässt sich die Hypothese wie folgt definieren:

$$h(\mathbf{x}_i) = \begin{cases} +1 & \text{wenn } \mathbf{w}^T \mathbf{x}_i + b \geq 0 \\ -1 & \text{wenn } \mathbf{w}^T \mathbf{x}_i + b \leq 0 \end{cases}$$

1.2 Minimierungsproblem

Die algorithmische Bestimmung der Parameter \mathbf{w} und b erfolgt über ein Minimierungsproblem. Die Spalte zwischen den Hilfsebenen H_+ und H_- soll maximiert werden. Über die Projektionseigenschaft des Skalarproduktes lässt sich die Breite der Spalte über den Ausdruck $\frac{2}{\|\mathbf{w}\|}$ bestimmen. Die Breite der Spalte soll nun maximiert werden, was zu folgendem Minimierungsproblem führt:

$$\begin{aligned} \max_{\mathbf{w}, b} \frac{2}{\|\mathbf{w}\|} &\Leftrightarrow \min_{\mathbf{w}, b} \|\mathbf{w}\| \Leftrightarrow \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{u.d.N. } &y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 \geq 0 \end{aligned}$$

Durch die Verwendung von Lagrange Multiplikatoren lässt sich das Minimierungsproblem in eine Form bringen, die nur noch von der Lagrangevariable α sowie den Supportvektoren bzw. den neu zu klassifizierenden Vektoren abhängt:

$$\begin{aligned} \mathbf{w} &= \sum_i \alpha_i y_i \mathbf{x}_i \\ \sum_i \alpha_i y_i &= 0 \\ h(\mathbf{x}) &= \begin{cases} +1 & \sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b \geq 0 \\ -1 & \sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b \leq 0 \end{cases} \end{aligned}$$

2 Kernel Trick

2.1 Abbildfunktion

Mit der gezeigten Definition von SVMs lassen sich linear separierbare Daten gut trennen. Es ist jedoch nicht möglich Daten mit komplexeren anordnungen mit hoher Wahrscheinlichkeit richtig zu klassifizieren.

Mit Hilfe einer Abbildfunktion $\phi(\vec{x})$ lässt sich diese einschränkung umgehen. Hierzu transformieren wir den Eingabevektor \vec{x} in einen neuen Raum in dem eine lineare Trennung möglich ist.

$$\begin{aligned} \phi : \mathbb{R}^n &\rightarrow \mathbb{R}^m \\ \mathbf{x} &\mapsto \mathbf{f} \end{aligned}$$

Um die transformierten Vektoren für die SVM zu nutzen ersetzen wir das Skalarprodukt in Minimierungs- und Hypothesenfunktion durch ein Skalarprodukt der transformierten Vektoren.

2.2 Kernel

Wenn die Dimension der transformierten Vektoren größer ist als die der Ursprungsvektoren erhöht dies den Rechenaufwand. Zusätzlich wird $\phi(\vec{x})$ nur für das Skalarprodukt genutzt.

Wir definieren daher einen so genannten Kernel, welcher das gleiche Ergebnis wie das Skalarprodukt der transformierten Vektoren aufweist, aber einfacher zu berechnen ist.

$$K(\mathbf{v}, \mathbf{w}) = \phi(\mathbf{v})^T \phi(\mathbf{w})$$

$$K : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$$

2.3 Verschiedene kernel

In der Praxis wird der lineare, der polynomielle und der Gauß'sche Kernel am meisten genutzt.

2.3.1 Linear

Der Lineare Kernel ist equivalent zu einer SVM ohne Kernel.

Er ist in der Lage linear trennbare Daten zu klassifizieren 2.

$$K(\mathbf{v}, \mathbf{w}) = \mathbf{v}^T \mathbf{w}$$

2.3.2 Polynomiell

Der polynomielle Kernel korrespondiert zu einem ϕ welches die Daten auf einer polynomiellen Fläche verteilt.

2.3.3 Gauß

Der Gauß'sche Kernel nutzt die Geometrische Distanz zu anderen Datenpunkten zur Klassifizierung. Um den maximalen Abstand der zu vergleichenden Punkte festzulegen kann der Parameter ϕ angepasst werden 5.

Er ermöglicht die Daten in verschiedenen Klustern zu trennen 4.

$$K(\mathbf{v}, \mathbf{w}) = \exp\left(-\frac{\|\mathbf{v} - \mathbf{w}\|^2}{2\sigma^2}\right)$$

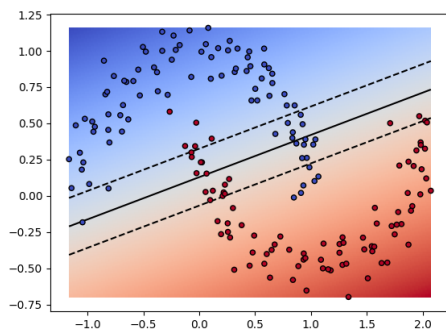


Abbildung 2: SVM mit linearem Kernel

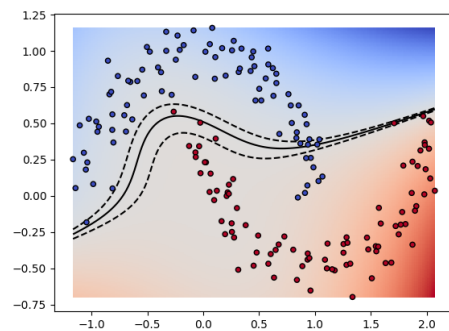


Abbildung 3: SVM mit Polynomiellen Kernel

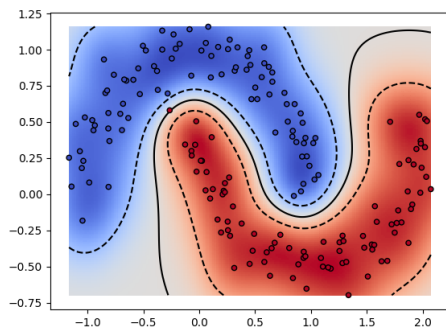


Abbildung 4: SVM mit Gauß'schem Kernel

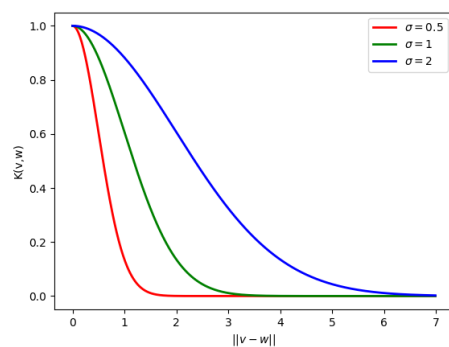


Abbildung 5: Verlauf des Gauß'schen Kerns

3 Tipps und Tricks zu Machine Learning und SVMs

Sämtlicher Code, Erläuterungen und das Präsentations-Skript sind online verfügbar: <https://github.com/NIPE-SYSTEMS/support-vector-machine>

Unten auf der Repository-Seite (siehe Link) befindet sich eine Beschreibung zu allen wichtigen Dateien im Repository.