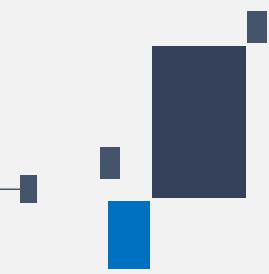




SVM原理阐述和实现

赵宇盛

西安电子科技大学
计算机科学与技术学院



SVM (Support Vector Machine) 支持向量机 ——用于监督学习的可定式的二分类器

通常有三种形式：

- 线性可分SVM
- 线性SVM
- 非线性SVM
- hard-margin SVM
- soft-margin SVM
- kernel SVM

SVM三个关键词

Margin

定义: γ 为整个数据集 D 中所有样本到分割超平面的最短距离 $\gamma = \min_n \gamma^{(n)}$

SVM的学习策略: **margin maximization**

\Leftrightarrow convex quadratic programming (with constraint)

\Leftrightarrow minimizing regularized hinge loss function (without constraint)

函数间隔和几何间隔

Dual

拉格朗日乘数法 优化问题 $(w, b) \rightarrow (\alpha)$

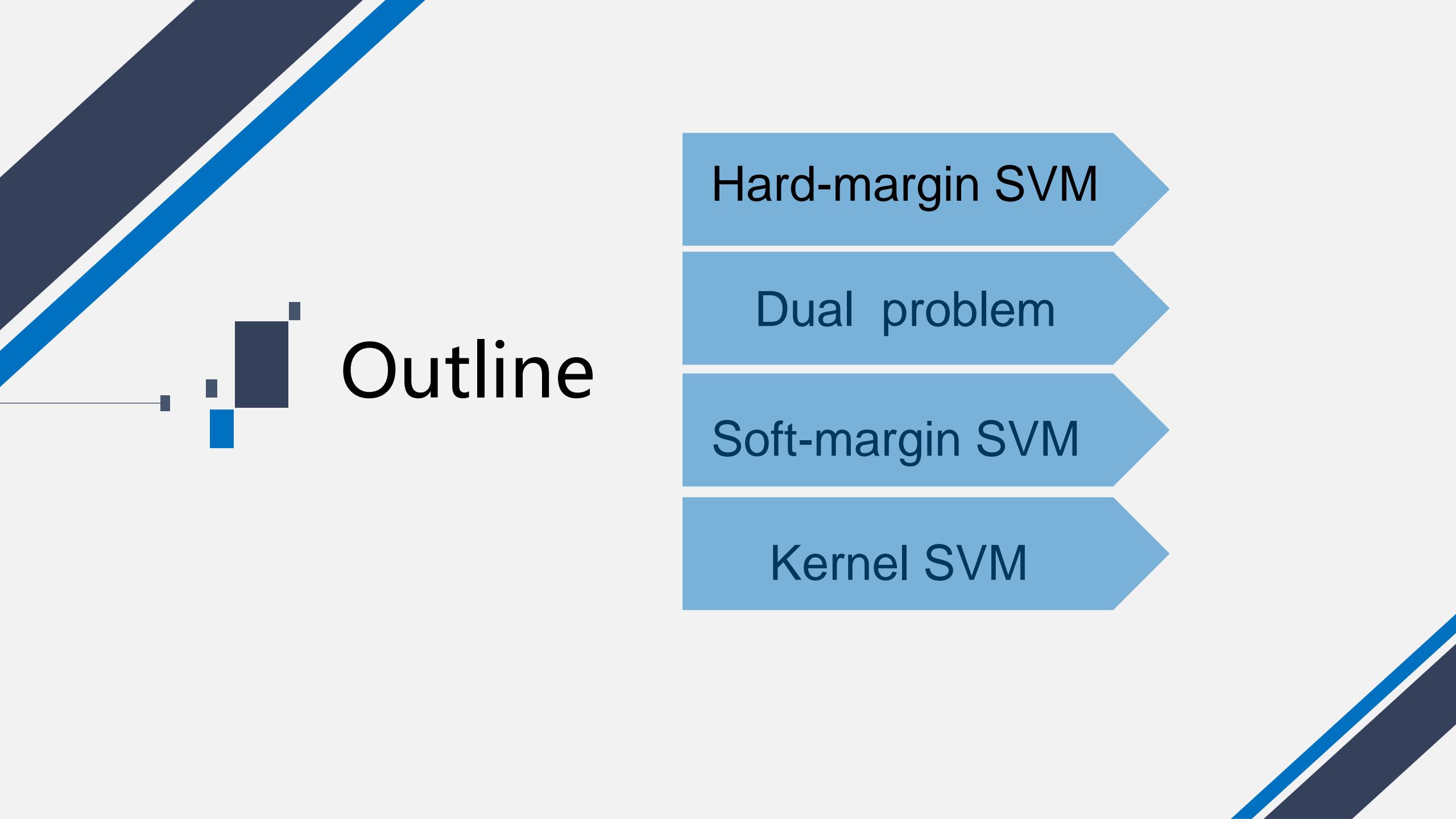
对偶问题消去了特征空间维度对求解超平面难度的影响

KKT条件中的互补松弛条件(complementary slackness)

Kernel Trick

让SVM在非线性可分场景得到适用

通过使用核函数表示特征向量之间的内积, 等价于隐式地在高维的特征空间中学习线性支持向量机。



Outline

Hard-margin SVM

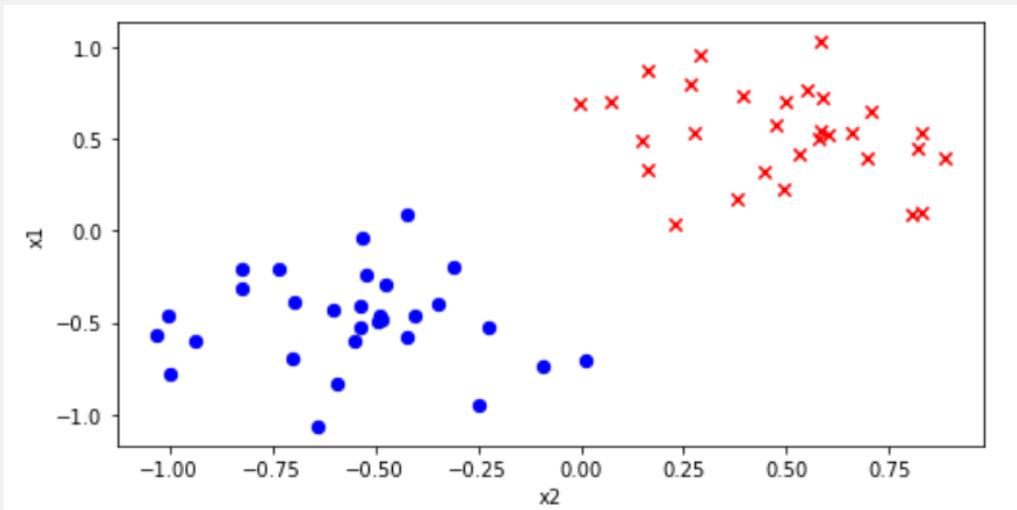
Dual problem

Soft-margin SVM

Kernel SVM

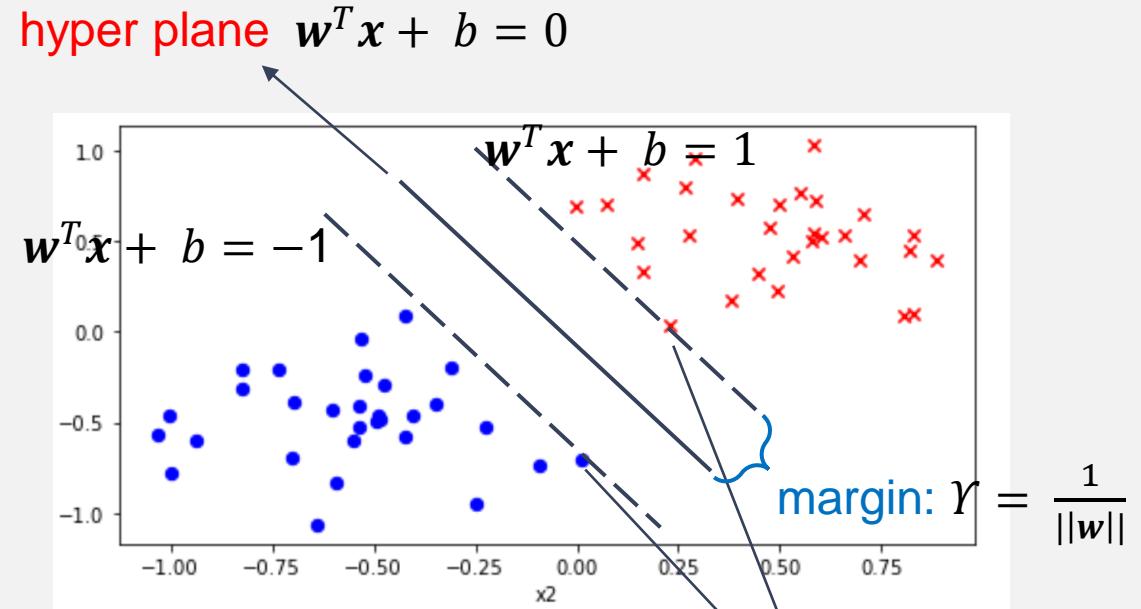
Hard-margin SVM

Example: 2-dim 线性可分情况



$$\{(x^i, y^i)\}_{i=0}^n, x^i \in \mathbb{R}^N, y^i \in \{-1, +1\}$$

上图 $N = 2$ 为例



SVM二分类器的**目标**:

找到某一个**hyper plane** (或者说某一组 w, b)

使得**margin**最大化

Hard-margin SVM

$$\max_{w,b} Y = \max_{w,b} \min_n Y^{(n)}$$

$$\begin{aligned} s.t. \quad & w^T x^i + b > 0, y^i = +1 \\ & w^T x^i + b < 0, y^i = -1 \end{aligned}$$

等价于

$$\begin{aligned} \max_{w,b} Y \equiv \max_w \frac{2}{||w||} \\ s.t. \quad & y^i(w^T x^i + b) > 0 \end{aligned}$$

为了保持所得到的hyper plane是唯一的 (准确地说, 参数组的唯一)

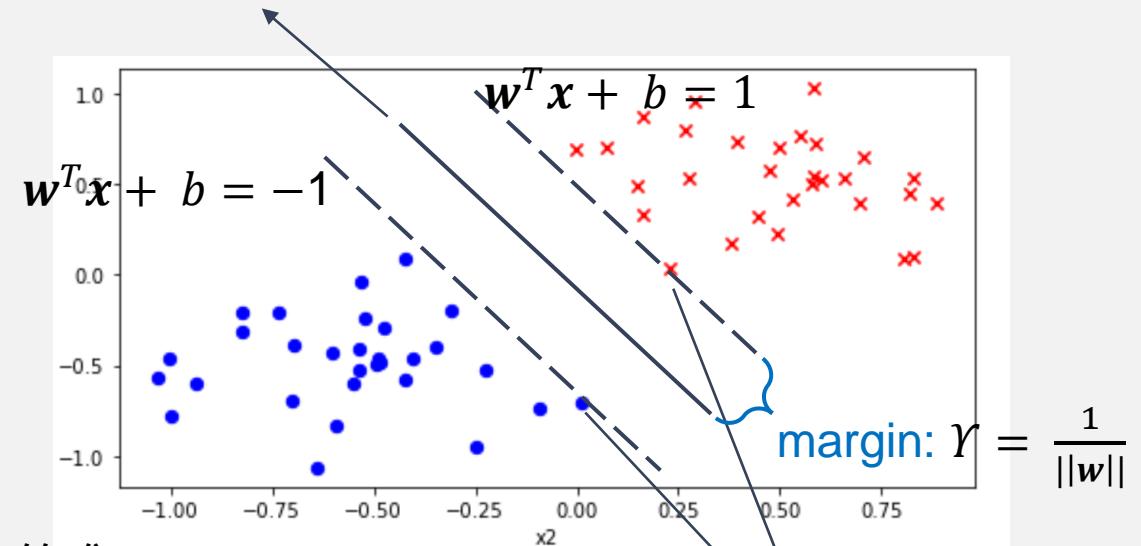
需要假设 $\min_w |w^T x^i + b| = 1$

规范化:

$$\frac{y^i(w^T x^i + b)}{||w||} \geq \frac{1}{||w||} = \gamma$$

因此, 以上等价于

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



support vector

$$\max_w \frac{2}{||w||^2}$$

s.t.

$$y^i(w^T x^i + b) \geq 1$$

当且仅当 (x^i, y^i) 为 support vector 时满足等价条件

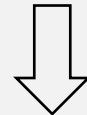
Hard-margin SVM

$$\begin{aligned} & \max_{\mathbf{w}} \frac{2}{\|\mathbf{w}\|^2} \\ \text{s.t.} \quad & y^i(\mathbf{w}^T \mathbf{x}^i + b) \geq 1 \end{aligned}$$

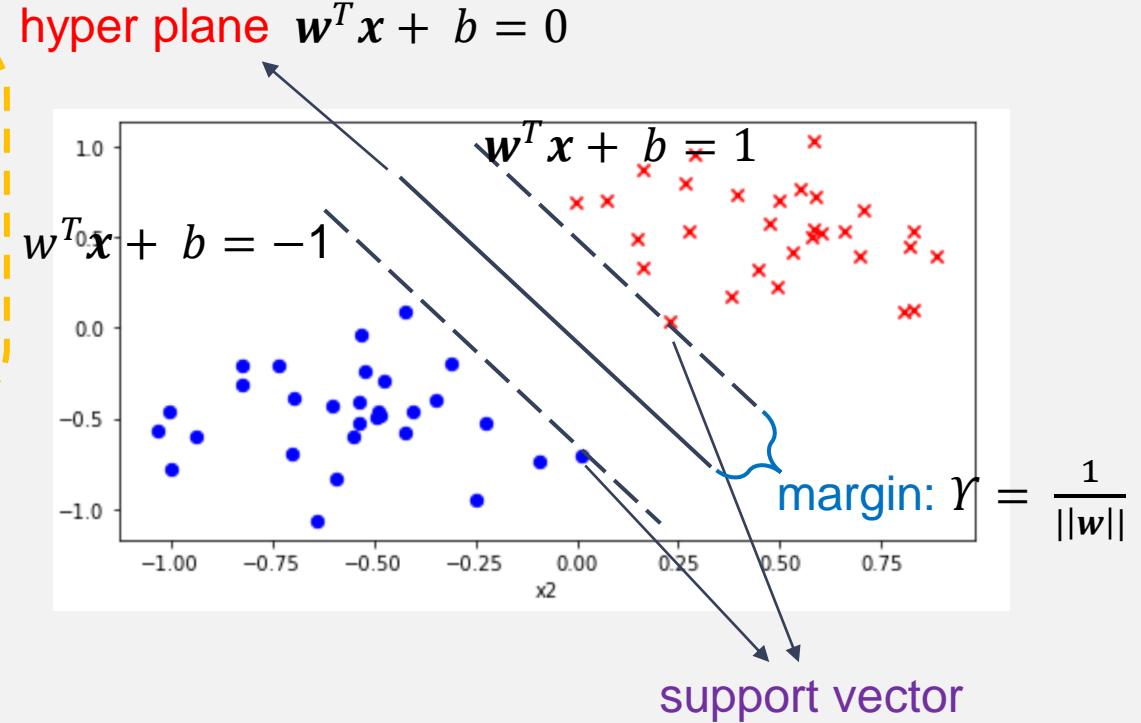
等价于

$$\begin{aligned} & \min_{\mathbf{w}} \frac{\|\mathbf{w}\|^2}{2} \\ \text{s.t.} \quad & y^i(\mathbf{w}^T \mathbf{x}^i + b) \geq 1 \end{aligned}$$

SVM: margin maximization



convex quadratic programming



如何去做这个凸优化问题呢？

Dual problem

不等式约束的优化问题

$$\begin{aligned} & \min_{\mathbf{w}} \frac{\|\mathbf{w}\|^2}{2} \\ & \text{s. t. } y^i(\mathbf{w}^T \mathbf{x}^i + b) \geq 1 \end{aligned}$$

拉格朗日乘数法

等式约束的优化问题

$$\begin{aligned} \Lambda(\mathbf{w}, b, \alpha) &= \frac{\|\mathbf{w}\|^2}{2} + \sum_{i=1}^n \alpha_i (1 - y^i(\mathbf{w}^T \mathbf{x}^i + b)) \\ \text{s. t. } \alpha_i &\geq 0, i = 1, 2, \dots, n \end{aligned}$$

分别求对 \mathbf{w}, b 的偏导

$$\begin{aligned} \frac{\partial \Lambda}{\partial \mathbf{w}} &= \mathbf{w} - \sum_{i=1}^n \alpha_i y^i \mathbf{x}^i & \mathbf{w} &= \sum_{i=1}^n \alpha_i y^i \mathbf{x}^i \\ \frac{\partial \Lambda}{\partial b} &= - \sum_{i=1}^n \alpha_i y^i & 0 &= \sum_{i=1}^n \alpha_i y^i \end{aligned}$$

代入 Λ 中得

$$\Gamma(\alpha) = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j ((\mathbf{x}^i)^T \mathbf{x}^j) + \sum_{i=1}^n \alpha_i \quad \text{s. t. } \sum_{i=1}^n \alpha_i y^i = 0$$



Dual problem

不等式约束的优化问题

$$\begin{aligned} & \min_{\mathbf{w}} \frac{\|\mathbf{w}\|^2}{2} \\ & \text{s. t. } y^i(\mathbf{w}^T \mathbf{x}^i + b) \geq 1 \end{aligned}$$

拉格朗日乘数法

等式约束的优化问题

$$\begin{aligned} & \Gamma(\alpha) = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j ((\mathbf{x}^i)^T \mathbf{x}^j) + \sum_{i=1}^n \alpha_i \\ & \text{s. t. } \sum_{i=1}^n \alpha_i y^i = 0, \quad \alpha_i \geq 0, i = 1, 2, \dots, n \end{aligned}$$

将原问题转换为对偶问题的变化：

原问题: 含有 $N + 1$ 个变量数, n 个约束条件。

对偶问题: 含有 n 个变量数, $n + 1$ 个约束条件。

N : \mathbf{w} 的维度, 即特征个数

n : 样本数量

当特征数远大于样本数($N \gg n$)时——拉格朗日对偶形式简化了原问题。

在约束条件下, 最大化该对偶函数依然是一个凸二次规划 (QP) 问题

Dual problem

不等式约束的优化问题

$$\begin{aligned} & \min_{\mathbf{w}} \frac{\|\mathbf{w}\|^2}{2} \\ & s.t. y^i(\mathbf{w}^T \mathbf{x}^i + b) \geq 1 \end{aligned}$$

拉格朗日乘数法

等式约束的优化问题

$$\begin{aligned} & \Gamma(\alpha) = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j ((\mathbf{x}^i)^T \mathbf{x}^j) + \sum_{i=1}^n \alpha_i \\ & s.t. \sum_{i=1}^n \alpha_i y^i = 0, \quad \alpha_i \geq 0, i = 1, 2, \dots, n \end{aligned}$$

拉格朗日对偶问题

KKT给出了互补松弛条件(complementary slackness)

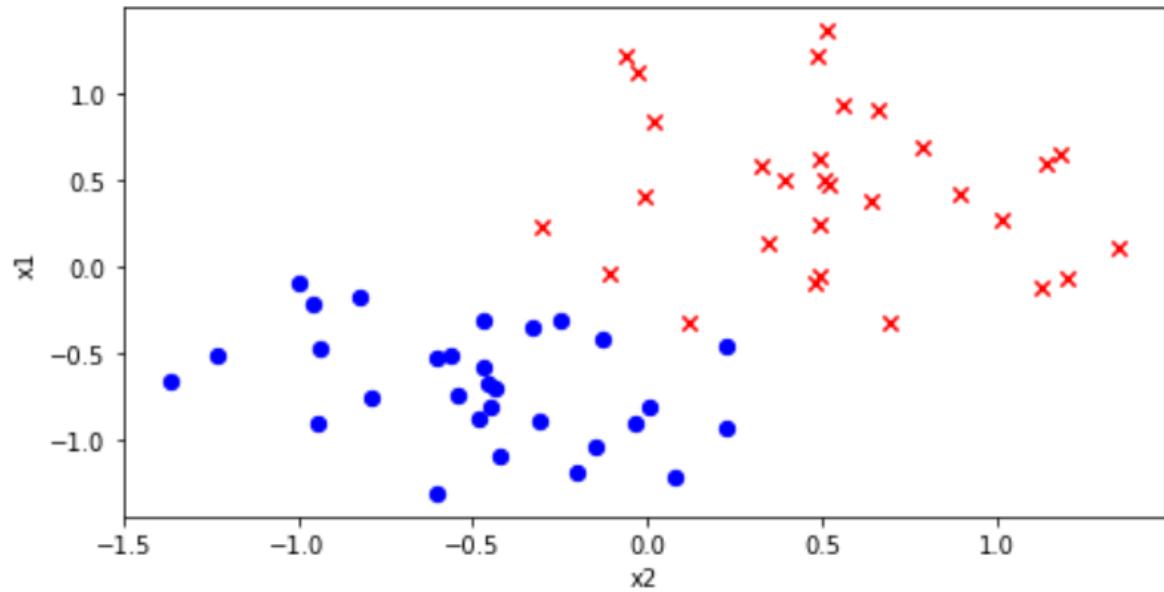
$$\alpha_i (1 - y^i (\mathbf{w}^T \mathbf{x}^i + b)) = 0, \quad \alpha_i \geq 0, i = 1, 2, \dots, n$$

互补松弛条件说明当最优解出现在不等式约束的内部，则约束失效。

当 $\alpha_i > 0$ 时, $y^i (\mathbf{w}^T \mathbf{x}^i + b) = 1$, (\mathbf{x}^i, y^i) 为决策边界上的support vector

对偶问题的最优解仅仅由support vector决定。

Soft-margin SVM



Hard-margin SVM

$$\begin{aligned} & \min_w \frac{\|w\|^2}{2} \\ & \text{s. t. } y^i(w^T x^i + b) \geq 1 \end{aligned}$$

优化问题调整为

正例和负例样本在特征空间中不是线性可分的情况

为了能够容忍部分不满足约束的样本，引入松弛变量 (Slack Variable) : ξ

Soft-margin SVM

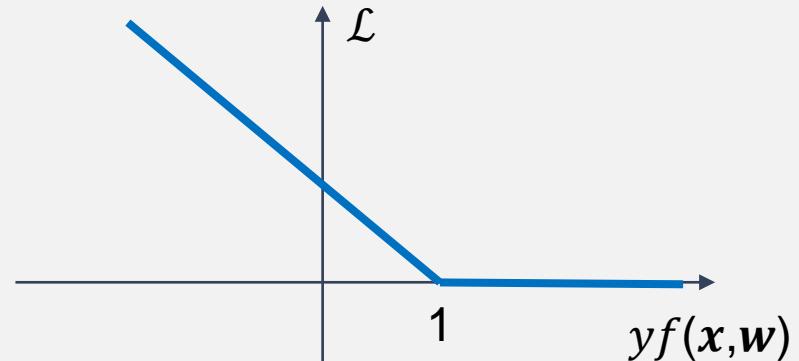
$$\begin{aligned} & \min_w \frac{\|w\|^2}{2} + C \sum_{i=1}^n \xi_i \\ & \text{s. t. } y^i(w^T x^i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0, i = 1, 2, \dots, n \end{aligned}$$

参数 $C > 0$ 用来控制间隔和松弛变量惩罚的平衡

Soft-margin SVM

引入hinge loss函数

$$\mathcal{L}_{hinge} = \max(0, 1 - y^i(\mathbf{w}^T \mathbf{x}^i + b)), i = 1, 2, \dots, n$$



Soft-margin SVM

表示为

$$\min_{\mathbf{w}} \frac{\|\mathbf{w}\|^2}{2} + C \sum_{i=1}^n \xi_i \quad \Leftrightarrow \quad \min_{\mathbf{w}} \sum_{i=1}^n \max(0, 1 - y^i(\mathbf{w}^T \mathbf{x}^i + b)) + \frac{\|\mathbf{w}\|^2}{2C}$$
$$\text{s.t. } y^i(\mathbf{w}^T \mathbf{x}^i + b) \geq 1 - \xi_i$$
$$\xi_i \geq 0, i = 1, 2, \dots, n$$

正则化项

从而，SVM的学习策略可以变为 **minimizing regularized hinge loss function**

无约束问题，可以用通用的深度学习方法来做！

例子：SVM的一个demo实现

基于

$$\min_w \sum_{i=1}^n \mathcal{L}_{hinge} + \left[\frac{1}{2C} \|\mathbf{w}\|^2 \right]$$

正则化项

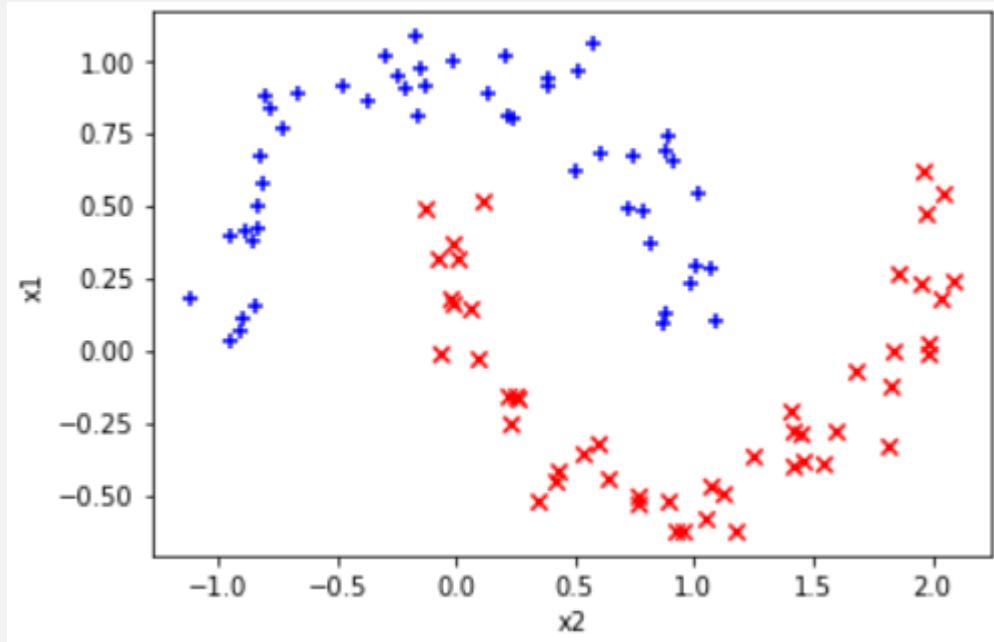


我的简单的小demo：（在线性可分的 \mathbb{R}^2 数据集上）

- Step 1: 未知参数 \mathbf{w}
- Step 2: 损失函数 \mathcal{L}_{hinge} + L2的正则项
- Step 3: 优化器: mini-batch SGD

Kernel SVM

非线性可分的情况，譬如



$$\Gamma(\alpha) = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j ((\mathbf{z}^i)^T \mathbf{z}^j) + \sum_{i=1}^n \alpha_i$$

$$\text{s. t. } \sum_{i=1}^n \alpha_i y^i = 0, \quad \alpha_i \geq 0, i = 1, 2, \dots, n$$

$$\Gamma(\alpha) = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j k(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i=1}^n \alpha_i$$

$$\text{s. t. } \sum_{i=1}^n \alpha_i y^i = 0, \quad \alpha_i \geq 0, i = 1, 2, \dots, n$$

使用核函数 (Kernel Function) 隐式地将样本从原始特征空间映射到更高维的空间

$$k(\mathbf{x}_i, \mathbf{x}_j) = (\Phi(\mathbf{x}_i))^T \Phi(\mathbf{x}_j)$$

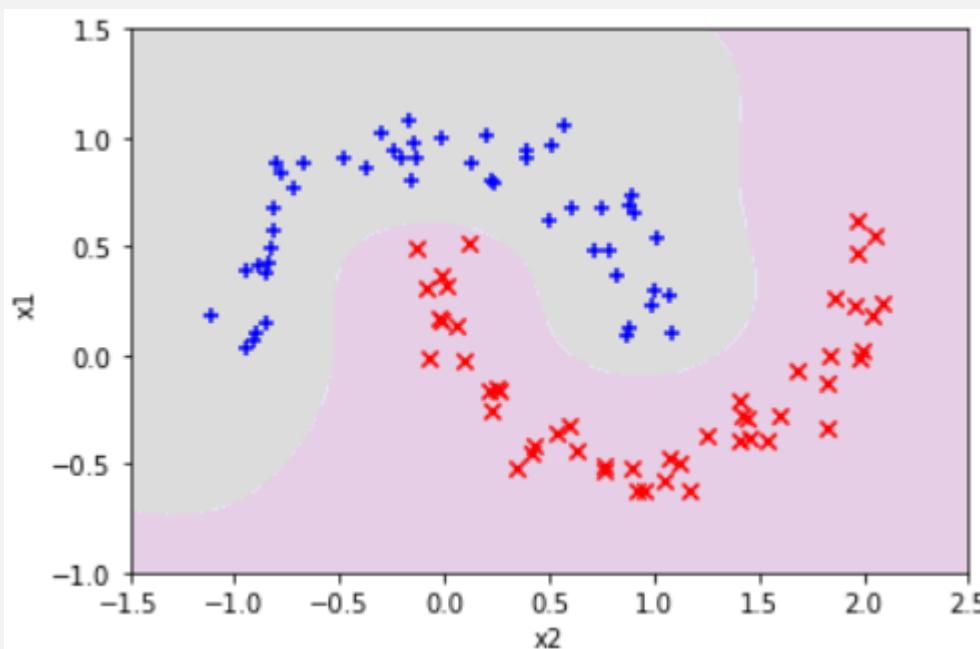
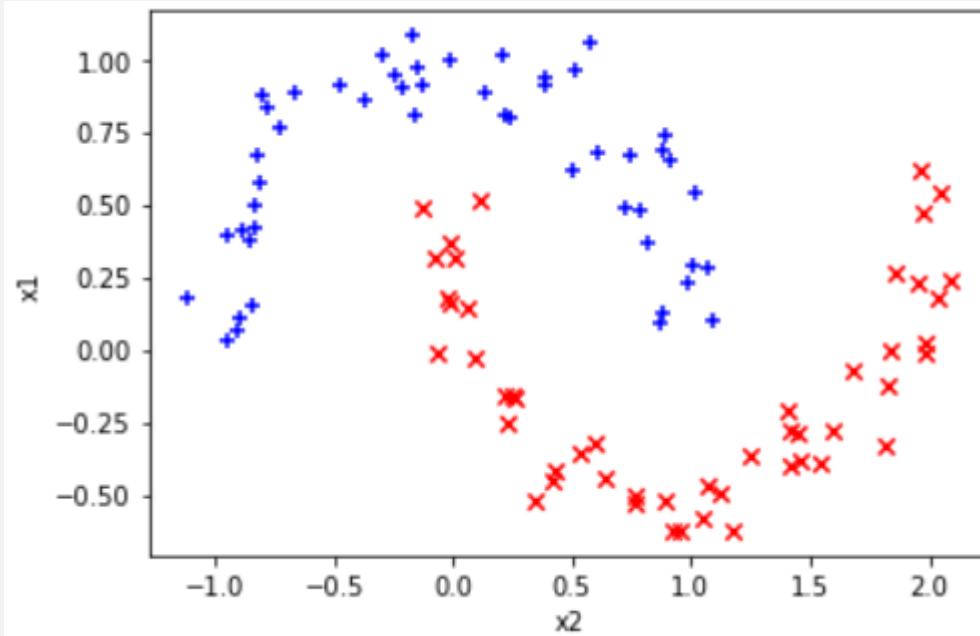
Kernel SVM

$$\begin{aligned}\Gamma(\alpha) = & -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j k(x_i, x_j) + \sum_{i=1}^n \alpha_i \\ \text{s. t. } & \sum_{i=1}^n \alpha_i y^i = 0, \quad \alpha_i \geq 0, i = 1, 2, \dots, n\end{aligned}$$

核函数 $k(x_i, x_j) = (\Phi(x_i))^T \Phi(x_j)$

不同的核函数其VC维也不同，对于 \mathbb{R}^N 大小的特征空间

- 对于线性核的分类器，其超平面是N-1维的，而VC维是N+1维
- 对于高斯核的分类器，其VC维是无穷





我的分享完毕

谢谢大家！