Machine Learning CS 4641



Support Vector Machine

Nakul Gopalan Georgia Tech

These slides are based on slides from Andrew Zisserman, Yaser S. Abu-Mostafa, Mahdi Roozbahani

Outline

- Precursor: Linear Classifier and Perceptron
- Support Vector Machine
- Parameter Learning

Binary Classification

Given training data (\mathbf{x}_i, y_i) for $i = 1 \dots N$, with $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \{-1, 1\}$, learn a classifier $f(\mathbf{x})$ such that

$$+ + f(\mathbf{x}_i) \begin{cases} \geq 0 & +1 \checkmark \\ < 0 & -1 \end{cases}$$

i.e. $y_i f(\mathbf{x}_i) > 0$ for a correct classification.



Linear Separability

linearly separable



not linearly separable





Linear Classifier



Linear Classifier (higher dimension)



Perceptron Algorithm

Input: A sequence of training examples $(\mathbf{x}_1, \mathbf{y}_1)$, $(\mathbf{x}_2, \mathbf{y}_2)$, \cdots where all $\mathbf{x}_i \in \Re^n$, $\mathbf{y}_i \in \{-1, 1\}$

- Initialize $\mathbf{w}_0 = \mathbf{0} \in \Re^n$
- For each training example (**x**_i, y_i):
 - Predict $\mathbf{y}' = \operatorname{sgn}(\mathbf{w}_t^T \mathbf{x}_i)$
 - If $y_i \neq y'$:
 - Update w_{t+1} ← w_t + r (y_i x_i)
- Return final weight vector

Mistake on positive: $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + r \mathbf{x}_i$ Mistake on negative: $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - r \mathbf{x}_i$

r is the learning rate, a small positive number less than 1

Update only on error. A mistakedriven algorithm

This is the simplest version. We will see more robust versions at the end

Mistake can be written as $y_i \mathbf{w}_t^T \mathbf{x}_i \leq 0$

Mistake on positive: $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + r \mathbf{x}_i$ Mistake on negative: $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - r \mathbf{x}_i$

Geometry of the perceptron update





All cases, error is zero and they are linear, so they are all good for generalization.

What is the Best θ ?



maximum margin solution: most stable under perturbations of the inputs



- if the data is linearly separable, then the algorithm will converge
- convergence can be slow ...
- separating line close to training data
- we would prefer a larger margin for generalization (better generalization)

Outline

- Precursor: Linear Classifier and Perceptron
- Support Vector Machine
- Parameter Learning



Finding θ with a **fat** margin



Computing the distance

The distance between x_i and the line $x\theta + b = 0$ where $|x_i\theta + b| = 1$ The vector θ is perpendicular to the decision line.

Consider x' and x'' on the plane

 $x'\theta + b = 0$ and $x''\theta + b = 0$ \downarrow $x'\theta + b = x''\theta + b$

$$\implies (x' - x'')\theta = 0$$



What is the distance of my fat margin?

What is the distance between x_i and the plane?

Let's take any point $\boldsymbol{\chi}$ on the line:

Distance would be projection of $(x_i - x)$ vector on θ .

To project the vector, we need to normalize θ to get the unit vector.

$$\hat{\theta} = \frac{\theta}{||\theta||} \Rightarrow \text{distance} = \left| (x_i - x)\hat{\theta} \right| \text{ which is the dot product}$$

$$x_2$$

$$x_1$$

$$x_i$$

What is the distance of my fat margin?

 $(\chi_{1}-\chi)\cdot \theta$

What is the distance between x_i and the plane?

Let's take any point $\boldsymbol{\chi}$ on the line:

Distance would be projection of $(x_i - x)$ vector on θ .

To project the vector, we need to normalize θ to get the unit vector.









 $\nabla (n) = \alpha \nabla g(x)$ $\frac{min}{g(x)} \begin{cases} (x) \\ g(x) \leq 0 \end{cases}$ Lagrange nul tipliers l'integrate $L(\chi, \chi, B) \ge f(\chi) + \alpha g(\chi) + \beta (h(\chi))$ $\frac{1}{2} (2, 0, 0, 0) = \frac{1}{2} (2, 0, 0) = \frac$ 7,5 $= \frac{100}{2} \cdot \frac{$

Krimal Problem: Optimal $d^{x} = max min \frac{100}{2} - \frac{N}{2} K; g(N;)$ $d^{x} = Max Min \frac{100}{2} - \frac{N}{2} K; g(N;)$ 501/ Dual problem: $L(0, b, X) = \frac{1}{2} \frac{00}{-2} = \frac{3}{2} \frac{1}{2} \frac{1$

Constrained optimization Minimize $\frac{1}{2} \theta \theta^T$ Subject to $y_i(x_i\theta + b) \ge 1$ for i = 1, 2, ..., N

 $\theta \in \mathbb{R}^d, b \in \mathbb{R}$

Using Lagrange method: But wait, there is an inequality in our constraints

We use Karush-Kuhn-Tucker (KKT) condition to deal with this problem

Constrained optimization Minimize $\frac{1}{2}\theta\theta^T$ Subject to $y_i(x_i\theta + b) \ge 1$ for i = 1, 2, ..., N $\theta \in \mathbb{R}^d, b \in \mathbb{R}$

Using Lagrange method: But wait, there is an inequality in our constraints

We use Karush-Kuhn-Tucker (KKT) condition to deal with this problem

$$g(x) = y_i(x_i\theta + b) - 1$$
 $\alpha = lagrange multiplier$

We need to optimize θ , b, and α

1) $g(x) \ge 0$ Primal feasibility

2) $\alpha \ge 0$ 3) $g(x)\alpha = 0$ Dual feasibility Complementary slackness $\Rightarrow \begin{cases} g(x) > 0, \quad \alpha = 0\\ \alpha > 0, \quad g(x) = 0 \end{cases}$



Lagrange formulation



$$\nabla_{\theta} \mathcal{L}(\theta, b, \alpha) = \theta - \sum_{i=1}^{N} \alpha_{i} y_{i} x_{i} = 0$$
$$\nabla_{b} \mathcal{L}(\theta, b, \alpha) = -\sum_{i=1}^{N} \alpha_{i} y_{i} = 0$$



Let's substitute these in the Lagrangian:

$$\mathcal{L}(\theta, b, \alpha) = \frac{1}{2}\theta\theta^{T} - \sum_{i=1}^{N} \alpha_{i}(y_{i}(x_{i}\theta + b) - 1)$$

$$\mathcal{L}(\theta, b, \alpha) = \sum_{i=1}^{N} \alpha_{i} + \frac{1}{2}\theta\theta^{T} - \sum_{i=1}^{N} \alpha_{i}(y_{i}(x_{i}\theta + b))^{T}$$

$$\mathcal{L}(\theta, b, \alpha) = \sum_{i=1}^{N} \alpha_{i} + \frac{1}{2}\theta\theta^{T} - \sum_{i=1}^{N} \alpha_{i}(y_{i}(x_{i}\theta)) = \sum_{i=1}^{N} \alpha_{i} + \frac{1}{2}\theta\theta^{T} - \theta\theta^{T} =$$

$$= \sum_{i=1}^{N} \alpha_{i} - \frac{1}{2}\theta\theta^{T}$$



The solution – quadratic programming

$$\max_{\alpha} \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} y_i y_j \alpha_i \alpha_j x_i x_j^T \qquad (1)$$

Quadratic programming packages usually use "min"





Training

$$\theta = \sum_{i=1}^{N} \frac{\alpha_i y_i x_i}{\sum_{i=1}^{N} \alpha_i y_i x_i}$$

No need to go over all datapoints

$$\rightarrow \theta = \sum_{x_i \text{ in } SV} \alpha_i y_i x_i$$

and for *b* pick any support vector and calculate: $y_i(x_i\theta + b) = 1$



Training

$$\theta = \sum_{i=1}^{N} \alpha_i y_i x_i$$

Testing

For a new test point s

Compute:

No need to go over all datapoints

$$\rightarrow \theta = \sum_{x_i \text{ in } SV} \alpha_i y_i x_i$$

and for *b* pick any support vector and calculate: $y_i(x_i\theta + b) = 1$

$$\mathbf{s}\boldsymbol{\theta} + \mathbf{b} = \sum_{x_i \text{ in } SV} \boldsymbol{\alpha}_i y_i x_i s^T + b$$

Classify s as class 1 if the result is positive, and class 2 otherwise

$$\sim$$

Geometric Interpretation





Regularization lecture

$$\theta = argmin_{\theta} E(\theta) = \frac{1}{n} \sum_{i=1}^{n} (y^{i} - z_{i}\theta)^{2} \longrightarrow MSE$$

Minimize
$$\frac{1}{N}(z\theta - y)^T(z\theta - y)$$

Subject to
$$\theta^t \theta \leq C$$

Subject to $\theta^t \theta \leq C$
Subject to $\theta^t \theta \leq C$
Subject to $\theta^t \theta \leq C$

For simplicity let's call θ_{lin} as weights' solution for non constrained one and θ for the constrained model.

Constrained optimization

Minimize
$$\frac{1}{2} \theta \theta^T$$

Subject to $y_i(x_i\theta + b) \ge 1$ for $i = 1, 2, ..., \theta \in \mathbb{R}^d, b \in \mathbb{R}$

Using Lagrange method: But wait, there is an inequality in our constraints

We use Karush-Kuhn-Tucker (KKT) condition to deal with this problem

Ν

Regularization vs SVMs

• Regularization:

$$Minimize \ \frac{1}{N}(z\theta - y)^T(z\theta - y)$$

Subject to $\theta^t \theta \leq C$

• SVM:

From x to z space



 (X_1^2, X_2^2)

In x space





$$\int_{0}^{\infty} points$$

let's say x is n × d
xx^T will be n × n

If I add millions of dimensions to x, would it affect the final size of xx^{T} ? In z space



 Z_1

In *x* space, they are called pre-images of support vectors



Out of sample error:

$$\mathbf{E}[E_{out}] \le \frac{\mathbf{E}[\# SV]}{N-1}$$

Take-Home Messages

- Linear Separability
- Perceptron
- SVM: Geometric Intuition and Formulation