

Linear SVM (continued)

from last time : $\vec{w}^T \vec{x}_\Delta + b \geq +1$ for all training data

 $\vec{w}^T \vec{x}_0 + b \leq -1$ for all training data
 $y = f(\Delta)$
 $y = -f(0)$

Trick : combine both constraints into one constraint

$y_i (\vec{w}^T \vec{x}_i + b) \geq 1$ for all training data
 "for all"

Email: tutorcs@163.com

Recall, for LINEAR classification, the boundary is that

QQ: 749389476

Hyperplane, ' H ', where $g(\vec{x}) = \vec{w}^T \vec{x} + b = 0 \quad \forall \vec{x} \in H$

<https://tutorcs.com>
 discriminant function \vec{w} define H this linear equation

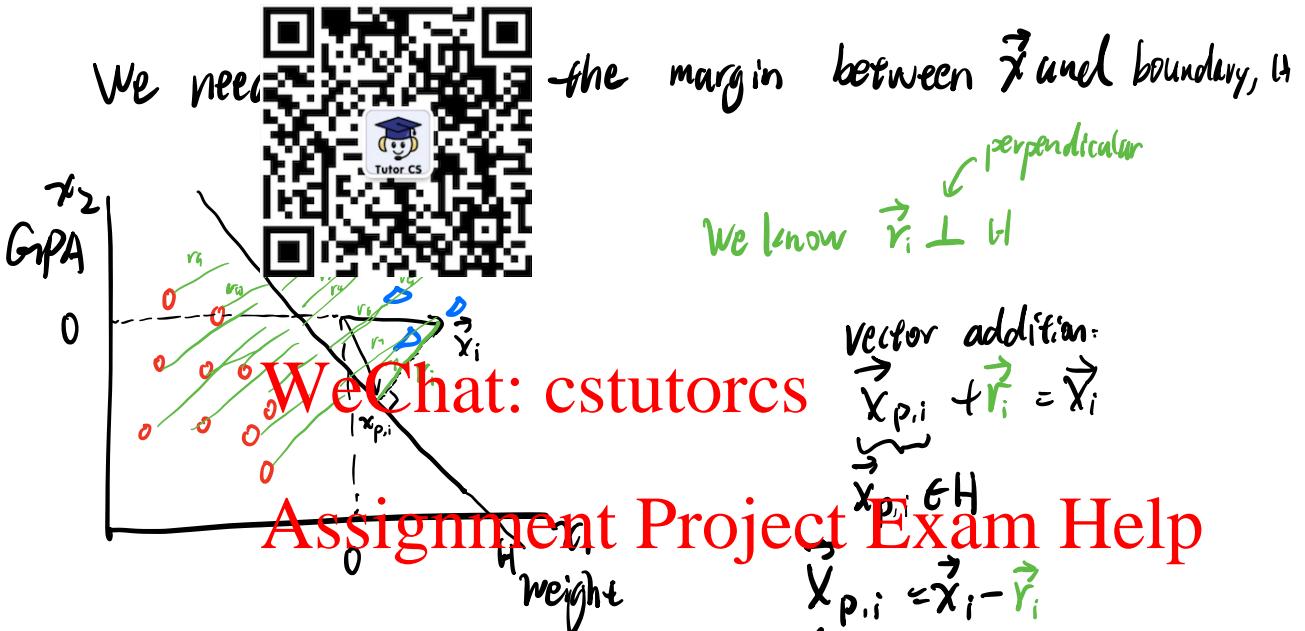
Say b is integrated into \vec{w}

$$\Rightarrow g_{\text{bias inside}}(\vec{x}) = \vec{w}_{\text{bias}} \vec{x} = 0 \quad \forall \vec{x} \in H$$

$$\hookrightarrow 2 \cdot \vec{w}_{\text{bias}} \vec{x} = 2 \cdot 0 = 0 \quad \text{still same boundary}$$

$$\Rightarrow 1000 \cdot \vec{w}_{\text{bias}} \vec{x} = 1000 \cdot 0 = 0$$

⇒ We can scale \vec{w} by any scalar value, and it still would define the same boundary, H .



Email: tutorcs@163.com

QQ: 749389476

$$\vec{w}^T \vec{x}_{p,i} + b = 0$$

$$\vec{w}^T (\vec{x}_{p,i} - \vec{r}_i) + b = 0$$

<https://tutorcs.com>

$$\vec{r}_i = \frac{\vec{r}_i}{\|\vec{w}\|} \vec{w}$$

we know from before
is our measure of confidence

$$\vec{w}^T \left(\vec{x} - \frac{\vec{r}_i}{\|\vec{w}\|} \vec{w} \right) + b = 0$$

We can find scalar distance r_i by solving this

Equation $\|\vec{w}\|^2$

$$\vec{w}^T \vec{x} - \frac{r_i}{\|\vec{w}\|} \vec{w}^T \vec{w} + b = 0$$

$$\Rightarrow \vec{w}^T \vec{x} - \frac{r_i}{\|\vec{w}\|} \|\vec{w}\|^2 + b = 0$$

$$\vec{w}^T \vec{x} + b = r_i \|\vec{w}\| \Rightarrow r_i = \frac{\vec{w}^T \vec{x} + b}{\|\vec{w}\|}$$

Def "margin for our SVM problem" = $\gamma(w, b) = \min_{i \in D} \gamma_i$

γ_i smallest of all y_i among our training data



Our job is to find H , (\vec{w} and b), that has the largest margin $\gamma(w, b)$ subject to our constrain (S)

goal : $\max_{\vec{w}, b} \gamma(w, b)$

Assignment Project Exam Help

Email: $tutorcs@163.com$

Subject to constraint QQ: 749389476

$$y_i: (\vec{w}^T \vec{x}_i + b) \geq 1 \quad \forall \vec{x}_i \in D$$

<https://tutorcs.com>

The same as the constrained optimization :

$$\max_{\vec{w}, b} \left\{ \frac{1}{\|\vec{w}\|} \min_{i \in D} |\vec{w}^T \vec{x}_i + b| \right\}$$

Subject to constrain

$$y_i: (\vec{w}^T \vec{x}_i + b) \geq 1 \quad \forall \vec{x}_i \in D$$

We know that we could solve it by any solver and it would yield the SAME "A":

⇒ Let's do it



This reduces our SVM optimization to
WeChat: cstutorcs

$\underset{\vec{w}, b}{\text{max}} \left\{ \frac{1}{\|\vec{w}\|^2} \right\}$ $\underset{\vec{w}, b}{\text{min}} \left\{ \frac{1}{2} \|\vec{w}\|^2 \right\}$

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

SVM formulation

<https://tutorcs.com>

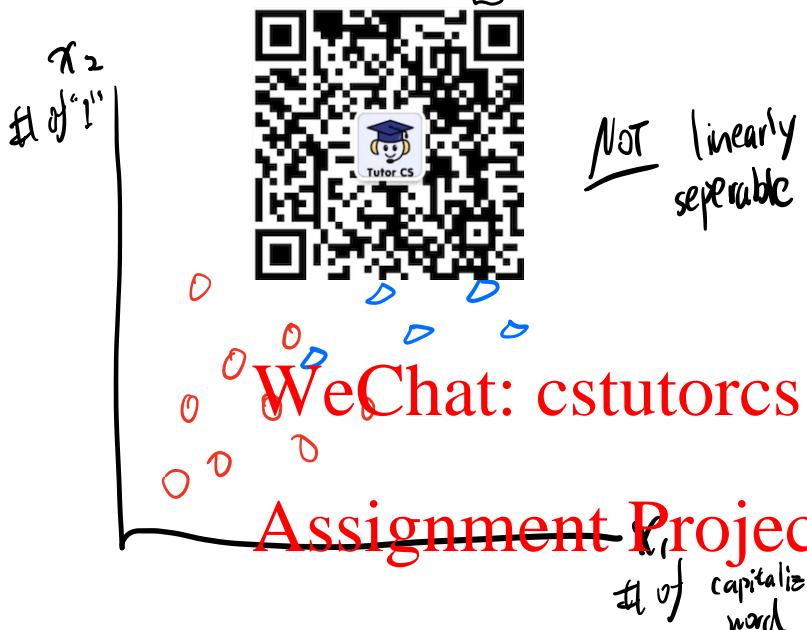
Linear SVM with slack

what if our training data is not

linearly separable?

Here, we define a cost function
that allows — but punishes — when
feature vectors fall too close to \vec{w}

(between the boundaries "gutters")
 程序代写与代做 CS 编程辅导
 and violate



SVM formulation with slack (non-differentiable)
 QQ: 749389476

$$(\vec{w}, b)_{\text{SVM}} = \arg \min_{\vec{w}, b} \left\{ \underbrace{\vec{w}^T \vec{w}}_{\text{regularizer}} + C \cdot \sum_{j=1}^N \max \{0, 1 - y_j (\vec{w}^T \vec{x}_j + b)\} \right\}$$

or equivalently (scale by $\frac{1}{C} = \lambda$)

we get the "Squared loss SVM" formulation (yes differentiable)

$$(\vec{w}, b)_{\text{SVM}} = \arg \min_{\vec{w}, b} \left\{ \lambda \vec{w}^T \vec{w} + \sum_{j=1}^N \left(\max \{0, 1 - y_j (\vec{w}^T \vec{x}_j + b)\} \right)^2 \right\}$$

Once we are equipped with those \vec{w} th b , we can find our decision boundary:

$$g(\vec{x}) = \vec{w}^T \vec{x} + b = 0 \quad \forall \vec{x} \in \mathcal{H}$$

For now (non-training) instances, we decide
程序代写代做 CS 编程辅导

$$g(x) \geq 0 \quad \text{and}$$



- "Hyperparameters" allows us to tune how important it is for us to allow violation of service constraints

- C WeChat: cstutorcs → you punishe the objective a lot

Assignment Project Exam Help
so minimizing ℓ_i is very important

- C Email: tutorcs@163.com
small → you don't care very much.....
where you know there are many misclassified instances in training data

<https://tutorcs.com>

In machine learning:

Def "Hyperparameter" = parameter that is not learned by the system, but is tuned to control the learning

Ex: C, λ , n
for SVM with slack
for regularization / learning rate for gradient descent

程序代写代做 CS 编程辅导



(0.001, 0.01, 0.1, 1, 10, 100, etc)

and then you test around the best of those

Assignment Project Exam Help

(0.1, 0.2, 0.3) etc.

Email: tutorcs@163.com

Non linear Binary Classification SVM: kernels

QQ: 749389476

Here: training data is not linearly separable:

<https://tutorcs.com>
and we are willing to relax linearly constraint

We want a (non-linear) transformation
to allow us to separate training instances
that are not linearly separable

Let's call the transformation that takes
each \vec{x} into this new higher dimensional space $h(\vec{x})$

We report SVM process based with transformed data, instead of raw data.

NOTE: I really care about $h(\vec{x})$... as long as some function (kernel)

which provides the dot product of these functions

WeChat: cstutorcs

DEF Kernel = measure of similarity used between insteads \vec{x} and some landmark \vec{l} where measure complies with Mercer's theorem

mercer's theorem: QQ: 749389476

a function $\phi(\vec{x}, \vec{l})$ is continuous, symmetric

($\phi(\vec{x}, \vec{l})$) is positive semi-definite

another space (possibly of higher dimension) such

that $\phi(\vec{x}, \vec{l}) = f_i(\vec{x})^T f_i(\vec{l})$

which kernels should we try first?

1) Linear kernel: no kernel: no transformation

2) Gaussian kernel = RBF (radial basis function)
 $= \phi(\vec{x}, \vec{l}) = e^{-\frac{\|\vec{x} - \vec{l}\|^2}{2\sigma^2}}$

3) Polynomial kernel $\phi(\vec{x}, \vec{l}) = (\vec{x}^\top \vec{l} + r)^d$

4) sigmoid kernel $\phi(\vec{x}, \vec{l}) = \tanh(\gamma \vec{x}^\top \vec{l} + r)$



$$\text{where } \tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$$

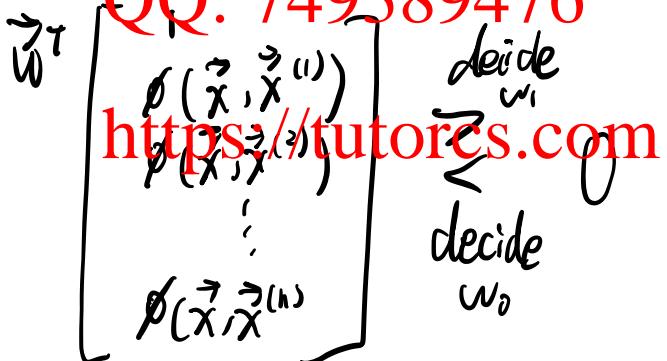
Note: If you use feature scaling, then similarity metrics will heavily bias in favor of features with large dynamic range.

→ You need ~~to apply feature scaling~~ Assignment Project Exam Help

landmarks you use are feature vectors from training data
Email: tutorcs@163.com

our new decision will be:

QQ: 749389476



Clustering

Unsupervised

no labels provided

in training data

Classification

not prediction

Algorithm

$D = \{ \vec{x}^{(i)} \}$ 程序代写代做CS编程辅导

Def class



Assign unlabeled data sample

to some similarity metric

k-means Algorithm

The most commonly used clustering algorithm.

WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>