

<https://powcoder.com>

Announcements

Assignment Project Exam Help

Reminder: ps4 self-grading form out, due Friday 10/30

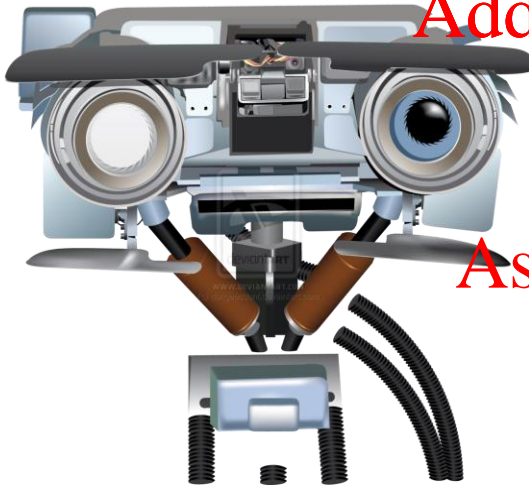
Assignment Project Exam Help

- pset 5 out today 10/29, due 11/5 (1 week)
- Midterm grades will go up by Monday (don't discuss it yet)

<https://powcoder.com>

Assignment Project Exam Help

Add WeChat powcoder



Assignment Project Exam Help

Support Vector Machines
<https://powcoder.com>

Add WeChat powcoder

CS542 Machine Learning

slides based on lecture by R. Urtasun

http://www.cs.toronto.edu/~urtasun/courses/CSC2515/CSC2515_Winter15.html

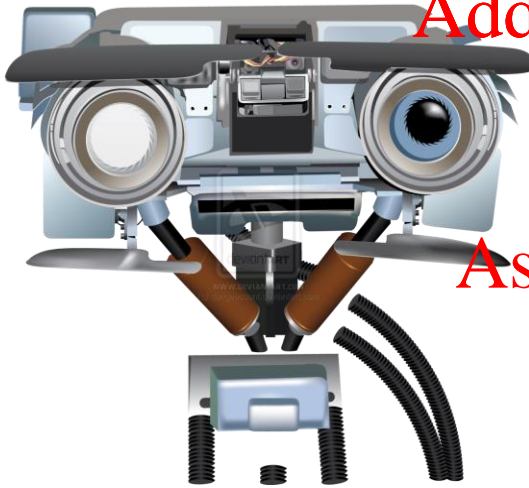
<https://powcoder.com>
Support Vector Machine (SVM)
Assignment Project Exam Help

- A *maximum margin* method, can be used for classification or regression
- SVMs can efficiently perform a non-linear classification using what is called the kernel trick, implicitly mapping their inputs into high-dimensional feature spaces
- First, we will derive *linear, hard-margin SVM* for linearly separable data, later for non-separable (soft margin SVM), and for nonlinear boundaries (kernel SVM)

<https://powcoder.com>

Assignment Project Exam Help

Add WeChat powcoder



Assignment Project Exam Help

Maximum Margin
<https://powcoder.com>

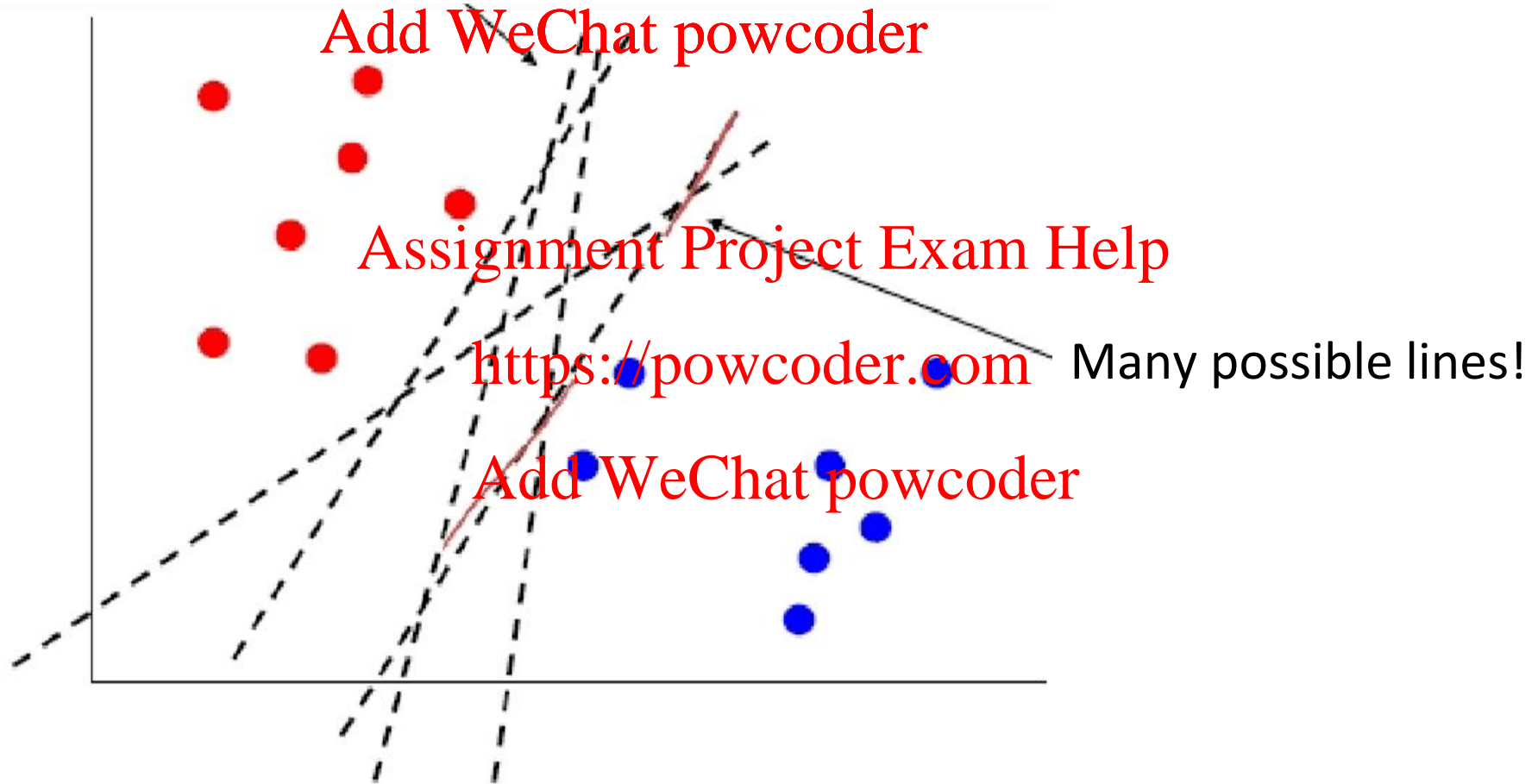
Add WeChat powcoder

<https://powcoder.com>
Recall: logistic regression

Assignment Project Exam Help

Decision boundary: $w^T x + b = 0$

Add WeChat powcoder



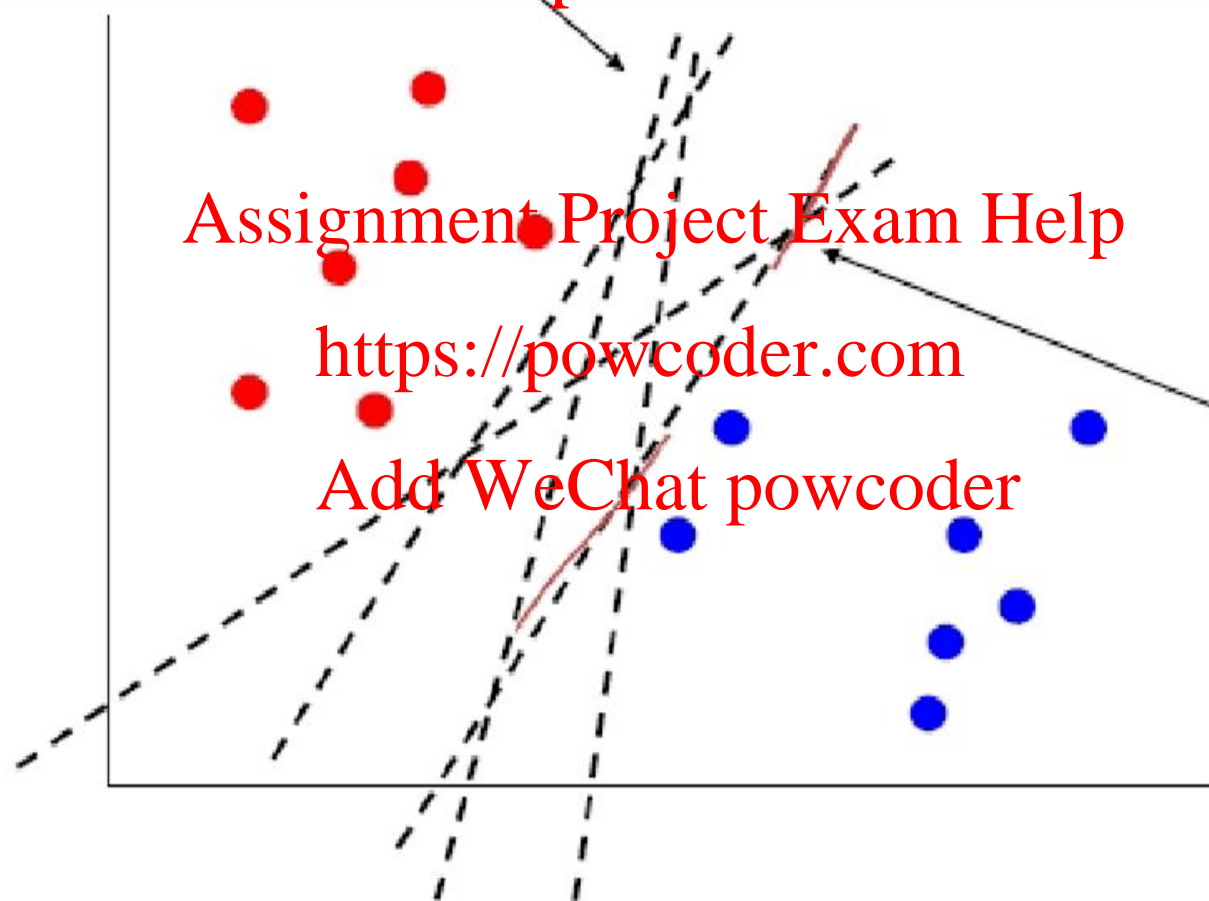
$$y = \begin{cases} +1 \text{ [red]} & \text{if } \text{sign}(\mathbf{w}^T \mathbf{x} + b) \geq 0 \\ -1 \text{ [blue]} & \text{if } \text{sign}(\mathbf{w}^T \mathbf{x} + b) < 0 \end{cases}$$

<https://powcoder.com>

Which classifier is best?

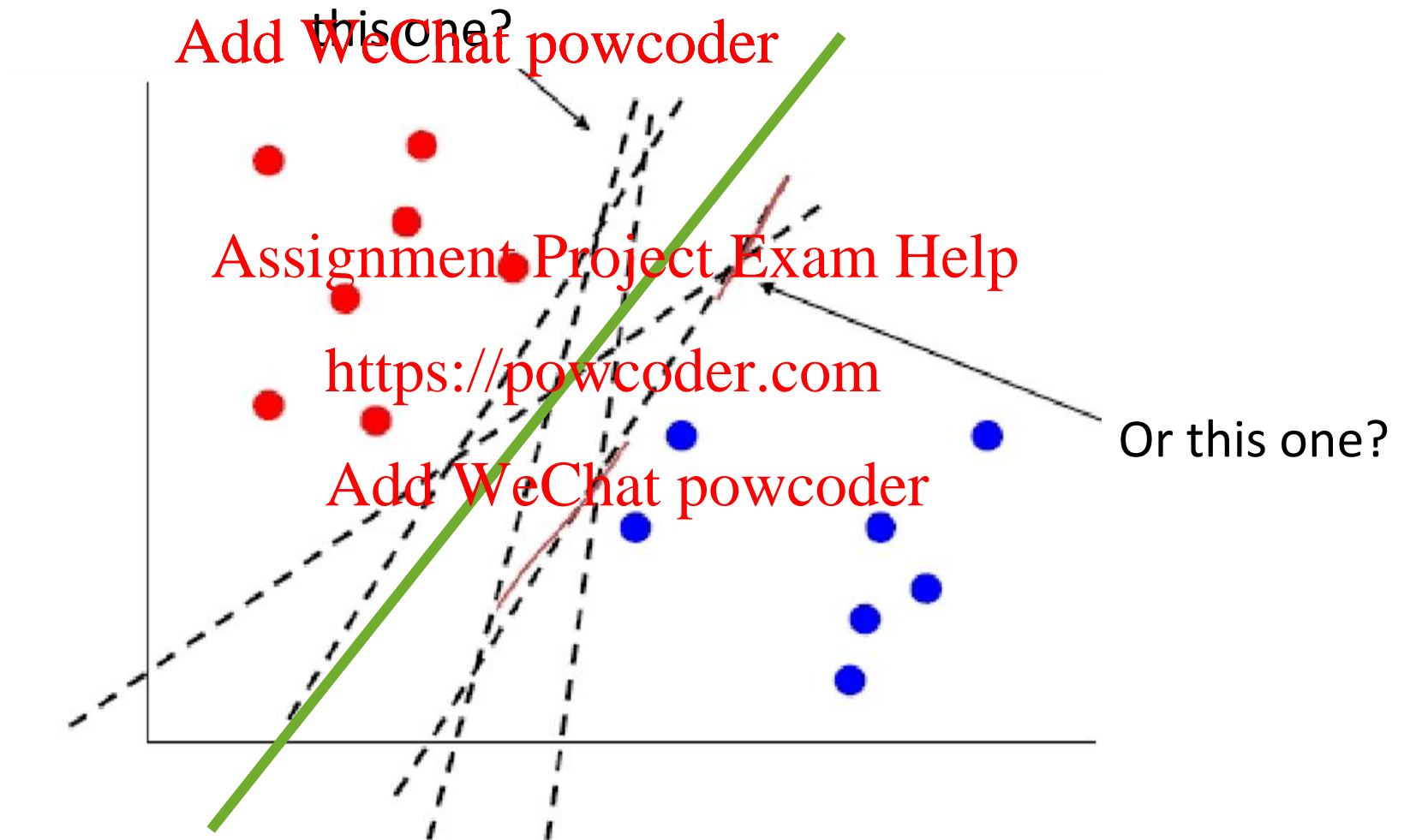
Assignment Project Exam Help

Add WeChat powcoder



Or this one?

<https://powcoder.com>
How about the one in the middle?
Assignment Project Exam Help



Intuitively, this classifier avoids misclassifying new test points generated from the same distribution as the training points

<https://powcoder.com> Max margin classification

Assignment Project Exam Help

Instead of fitting all the points, focus on boundary points

Add WeChat powcoder

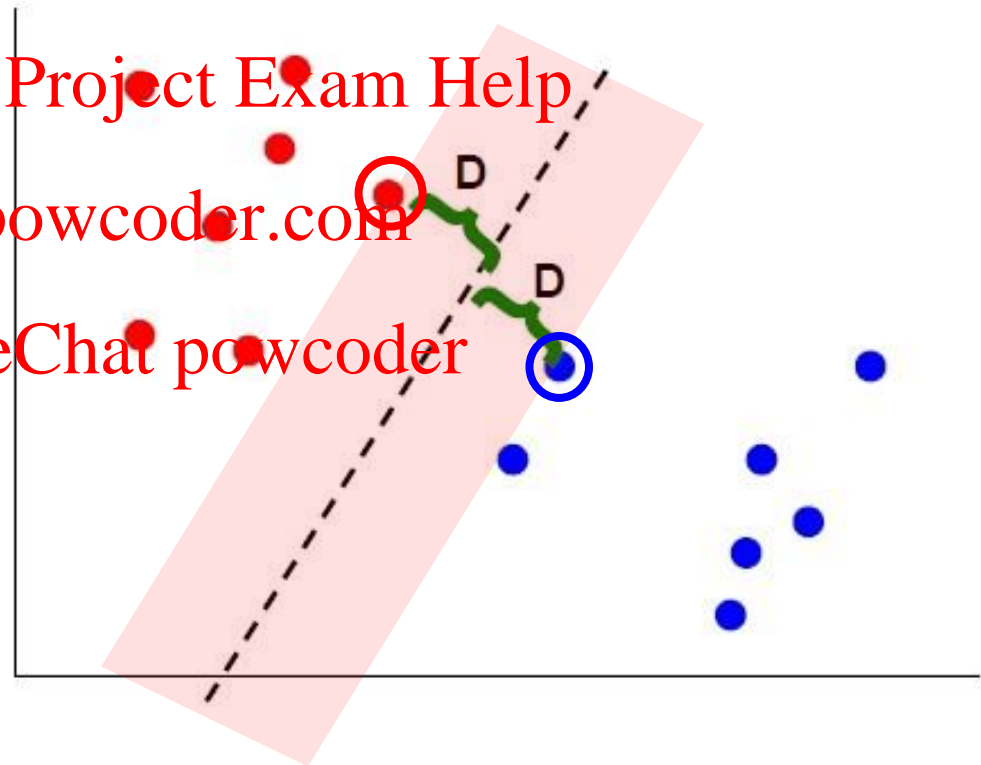
Aim: learn a boundary that leads to the largest margin (buffer)
from points on both sides

Why: intuition; theoretical

support: robust to small

perturbations near the
boundary

And works well in practice!

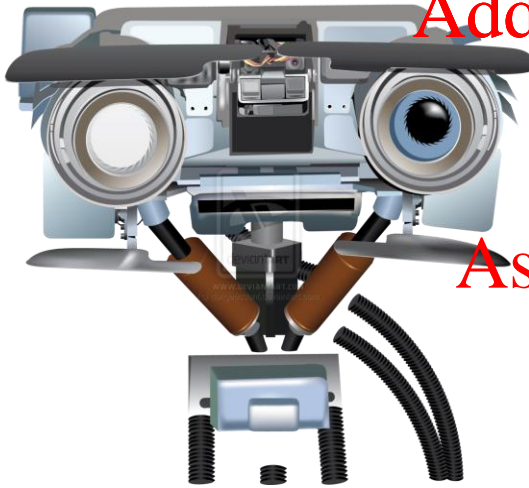


Subset of vectors that support (determine boundary) are called the
support vectors (circled)

<https://powcoder.com>

Assignment Project Exam Help

Add WeChat powcoder



Assignment Project Exam Help

Max-Margin Classifier
<https://powcoder.com>

Add WeChat powcoder

<https://powcoder.com>

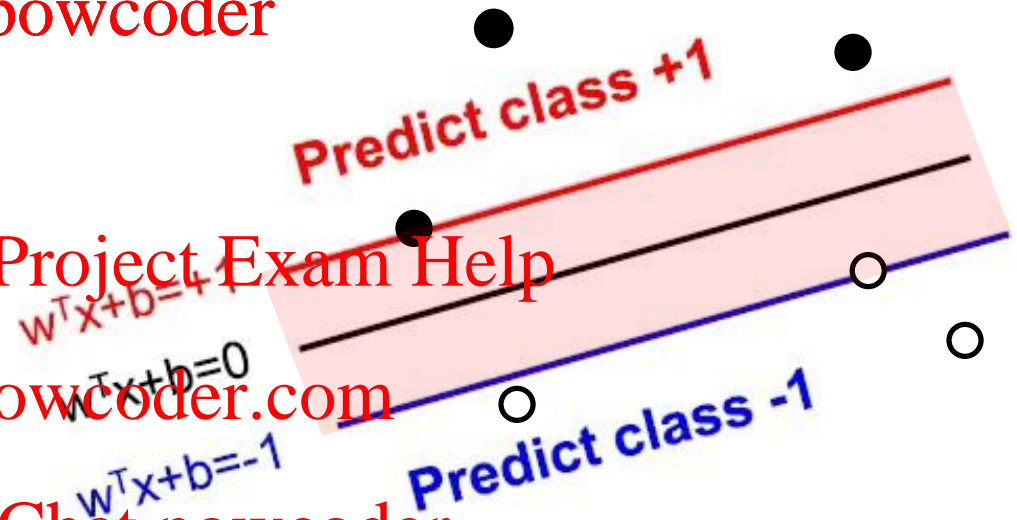
Max Margin Classifier

Assignment Project Exam Help

“Expand” the decision boundary to include a margin (until we hit first point on either side)

Use margin of 1

Inputs in the margins are of unknown class



Classify as +1

if

$$w^T x + b \geq 1$$

Classify as -1

if

$$w^T x + b \leq -1$$

Undefined

if

$$-1 < w^T x + b < 1$$

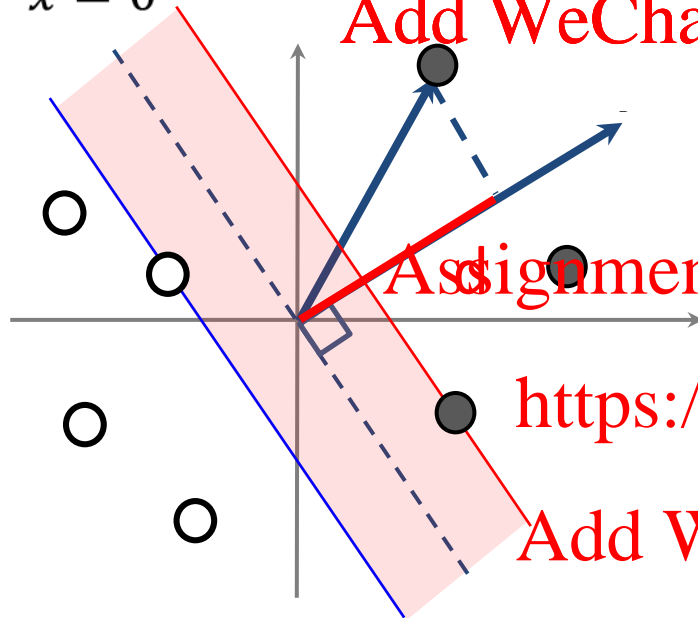
<https://powcoder.com>
Why is the margin = 1?

Assignment Project Exam Help

Decision boundary

$$w^T x = 0$$

Add WeChat powcoder



- Assume $b = 0$ for simplicity
- w is orthogonal to the decision plane
- Scaling margin and weight vector by the same constant $c > 0$ does not change the inequality

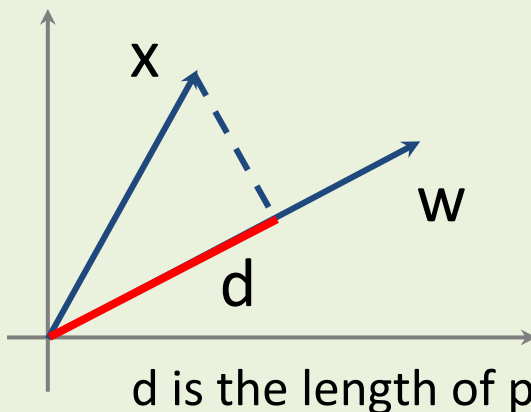
Assignment Project Exam Help

<https://powcoder.com>

$$w^T x \geq 1$$

$$c * w^T x \geq 1 * c$$

Add WeChat powcoder



Aside: vector inner product

$$\begin{aligned} w^T x &= d \|w\|_2 = \\ &= w_1 x_1 + w_2 x_2 \end{aligned}$$

$$d = \frac{w^T x}{\|w\|_2}$$

<https://powcoder.com>

Computing the Margin

Assignment Project Exam Help

First note that the w vector is orthogonal to the +1 plane

If u and v are two points on that plane, then $w^T(u-v) = 0$

Same is true for -1 plane

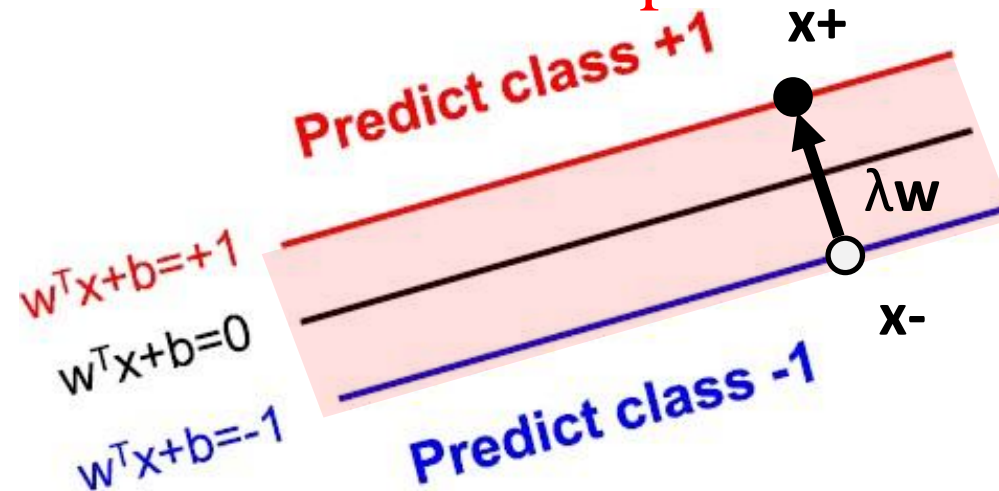
Assignment Project Exam Help

Also: for point x^+ on +1 plane and x^- nearest point on -1 plane:

$$x^+ = \lambda w + x^-$$

<https://powcoder.com>

Add WeChat powcoder



<https://powcoder.com>

Computing the Margin

Assignment Project Exam Help

Also: for point \mathbf{x}^+ on +1 plane and \mathbf{x}^- nearest point on -1 plane:

$$\mathbf{x}^+ = \lambda \mathbf{w} + \mathbf{x}^-$$

$$\mathbf{w}^T \mathbf{x}^+ + b = 1$$

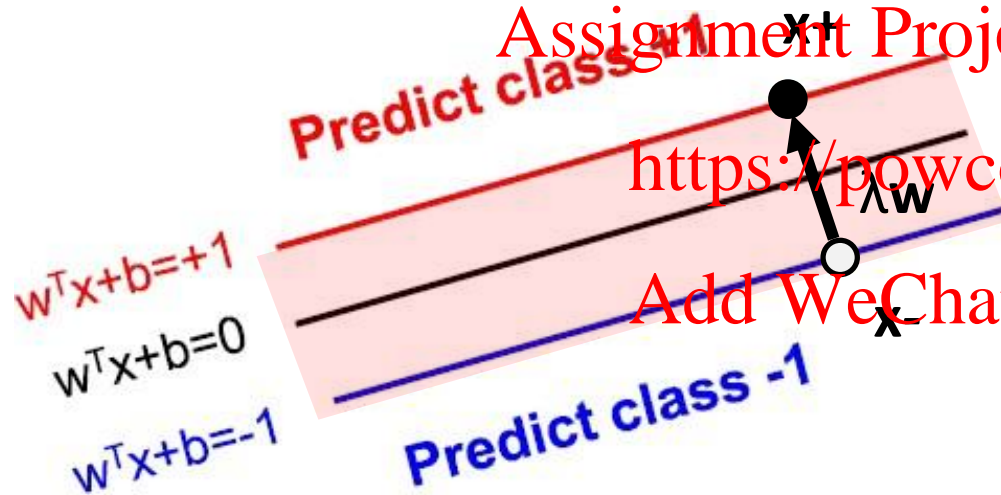
$$\mathbf{w}^T (\lambda \mathbf{w} + \mathbf{x}^-) + b = 1$$

$$\mathbf{w}^T \mathbf{x}^- + b + \lambda \mathbf{w}^T \mathbf{w} = 1$$

$$-1 + \lambda \mathbf{w}^T \mathbf{w} = 1$$

$$\lambda = \frac{2}{\mathbf{w}^T \mathbf{w}}$$

→ inversely proportional to $\mathbf{w}^T \mathbf{w}$, the square of the length of \mathbf{w}



<https://powcoder.com>

Computing the Margin

Assignment Project Exam Help

Define the margin M to be the distance between the +1 and -1 planes

Add WeChat powcoder

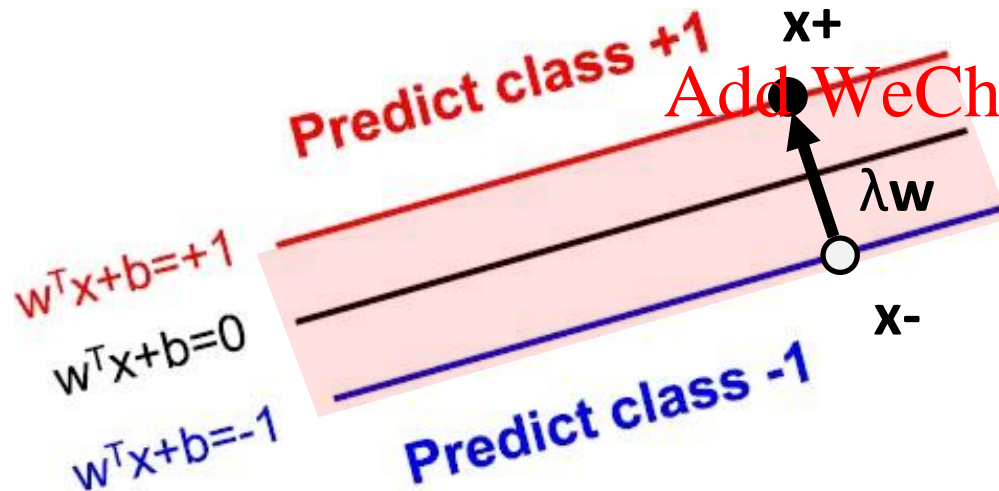
We can now express this in terms of \mathbf{w} □

to maximize the margin we minimize the length of \mathbf{w}

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



$$M = \|\mathbf{x}^+ - \mathbf{x}^-\|$$

$$= \|\lambda \mathbf{w}\| = \lambda \sqrt{\mathbf{w}^T \mathbf{w}}$$

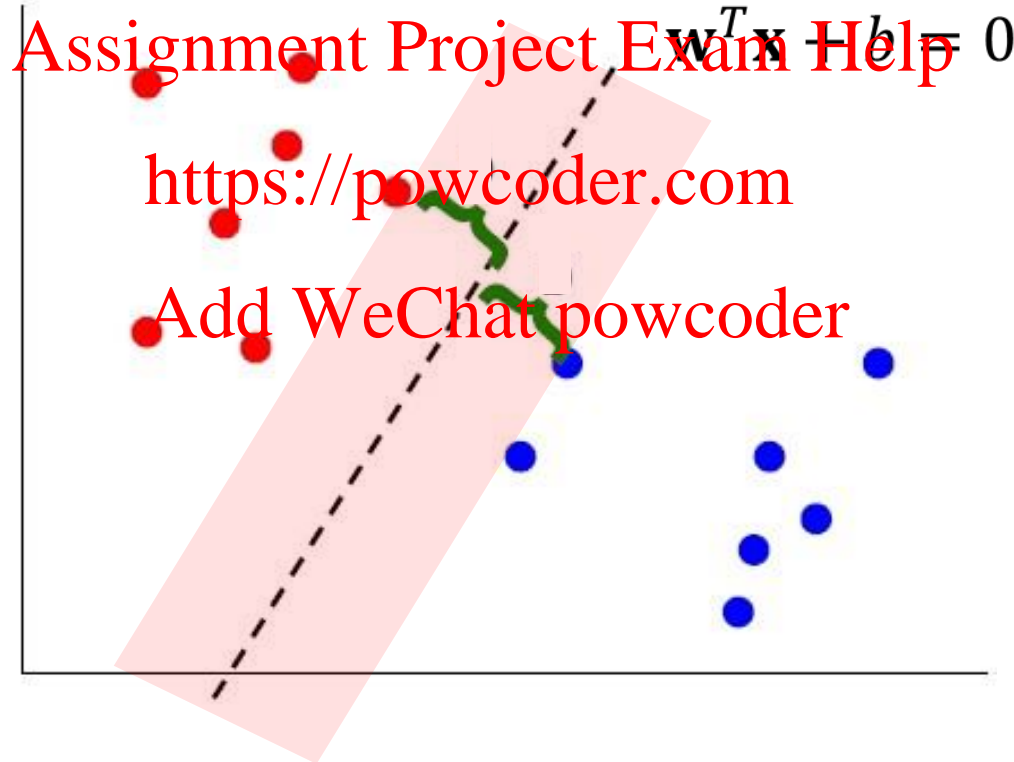
$$= 2 \frac{\sqrt{\mathbf{w}^T \mathbf{w}}}{\mathbf{w}^T \mathbf{w}} = \frac{2}{\sqrt{\mathbf{w}^T \mathbf{w}}}$$

<https://powcoder.com>
Maximizing the margin is equivalent to
regularization

Assignment Project Exam Help

Add WeChat powcoder

To maximize the margin we minimize the length of \mathbf{w} , or $\|\mathbf{w}\|^2$

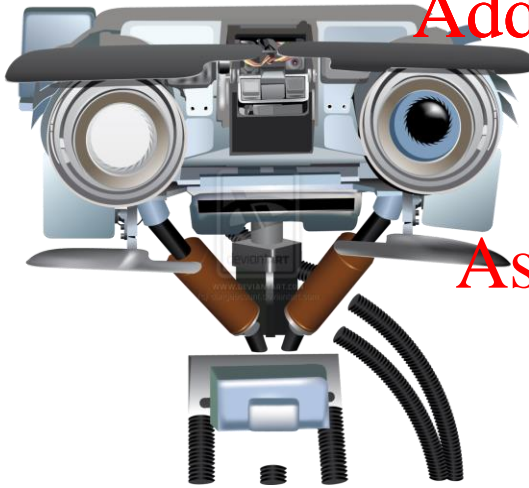


But not same as regularized logistic regression, the SVM loss is different! Only care about boundary points.

<https://powcoder.com>

Assignment Project Exam Help

Add WeChat powcoder



Assignment Project Exam Help

Linear SVM

<https://powcoder.com>

Add WeChat powcoder

<https://powcoder.com>

Linear SVM Formulation

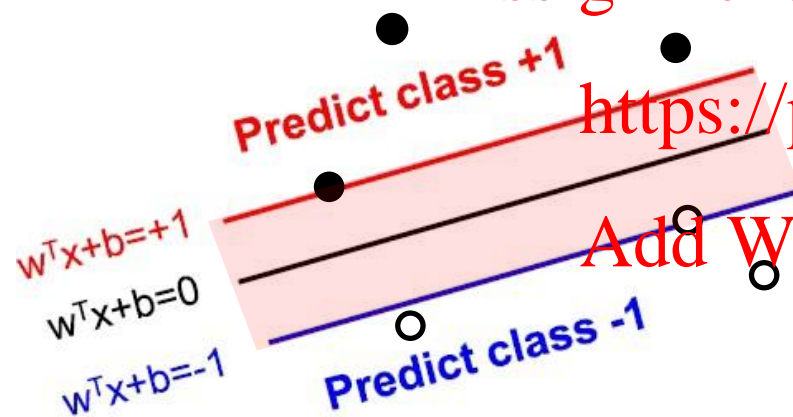
Assignment Project Exam Help

We can search for the optimal parameters (\mathbf{w} and b) by finding a solution that:

[Add WeChat powcoder](https://powcoder.com)

1. Correctly classifies the training examples: $\{x_i, y_i\}, i=1, \dots, n$
2. Maximizes the margin (same as minimizing $\|\mathbf{w}\|^2$)

Assignment Project Exam Help



<https://powcoder.com>

[Add WeChat powcoder](https://powcoder.com)

$$\min \frac{1}{2} \|\mathbf{w}\|^2$$

$$s.t. (\mathbf{w}^T \mathbf{x}_i + b) y_i \geq 1 \quad \forall i$$

This is the **primal formulation**, can be optimized via gradient descent, EM, etc.

Apply Lagrange multipliers: formulate equivalent problem

<https://powcoder.com>

Lagrange Multipliers

Assignment Project Exam Help

Convert the primal constrained minimization to an unconstrained optimization problem: represent constraints as penalty terms:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + \textit{penalty_term}$$

Assignment Project Exam Help

For data $\{(\mathbf{x}_i, y_i)\}$ use the following penalty term:

$$\begin{cases} 0 & \text{if } (\mathbf{w}^T \mathbf{x}_i + b)y_i \geq 1 \\ \infty & \text{otherwise} \end{cases} = \max_{\alpha_i \geq 0} \alpha_i [1 - (\mathbf{w}^T \mathbf{x}_i + b)y_i]$$

<https://powcoder.com>
Add WeChat powcoder

≤ 0 if constraint satisfied

Introduced Lagrange variables $\alpha_i \geq 0$; find ones that maximize term:

- If a constraint is satisfied, large α_i ensures smaller penalty
- If a constraint is violated, large α_i ensures larger penalty

Note, we are now minimizing with respect to \mathbf{w} and b , and maximizing with respect to \mathbf{a} (additional parameters)

<https://powcoder.com>

Lagrange Multipliers

Assignment Project Exam Help

Convert the primal constrained minimization to an unconstrained optimization problem: represent constraints as penalty terms:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + \textit{penalty_term}$$

Assignment Project Exam Help

For data $\{(\mathbf{x}_i, y_i)\}$ use the following penalty term:

$$\begin{cases} 0 & \text{if } (\mathbf{w}^T \mathbf{x}_i + b)y_i \geq 1 \\ \infty & \text{otherwise} \end{cases} = \max_{\alpha_i \geq 0} \alpha_i [1 - (\mathbf{w}^T \mathbf{x}_i + b)y_i]$$

Rewrite the minimization problem:

$$\min_{\mathbf{w}, b} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n \max_{\alpha_i \geq 0} \alpha_i [1 - (\mathbf{w}^T \mathbf{x}_i + b)y_i] \right\}$$

Where $\{\alpha_i\}$ are the

Lagrange multipliers

$$\min_{\mathbf{w}, b} \max_{\alpha_i \geq 0} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n \alpha_i [1 - (\mathbf{w}^T \mathbf{x}_i + b)y_i] \right\}$$

<https://powcoder.com>
Solution to Linear SVM
 Assignment Project Exam Help

Swap the 'max' and 'min':

$$\max_{\alpha_i \geq 0} \min_{\mathbf{w}, b} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i [1 - (\mathbf{w}^T \mathbf{x}_i + b) y_i] \right\}$$

$$= \max_{\alpha \in \mathbb{R}^n} \min_{\mathbf{w}, b} J(\mathbf{w}, b; \alpha)$$

First minimize $J()$ w.r.t. $\{\mathbf{w}, b\}$ for any fixed setting of the Lagrange multipliers:

$$\frac{\partial}{\partial \mathbf{w}} J(\mathbf{w}, b; \alpha) = \mathbf{w} - \sum_{i=1}^n \alpha_i \mathbf{x}_i y_i = 0$$

$$\frac{\partial}{\partial b} J(\mathbf{w}, b; \alpha) = - \sum_{i=1}^n \alpha_i y_i = 0$$

Then substitute back into $J()$ and simplify to get final optimization:

$$L = \max_{\alpha_i \geq 0} \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j (\mathbf{x}_i \cdot \mathbf{x}_j) \right\}$$

<https://powcoder.com>

Dual Problem

Assignment Project Exam Help

Final optimization: maximize this loss over α_i 's: only dot products of

data points needed

$$L = \max_{\alpha_i \geq 0} \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j (\mathbf{x}_i \cdot \mathbf{x}_j) \right\}$$

Assignment Project Exam Help

<https://powcoder.com>

$$\text{subject to } \alpha_i \geq 0, \quad \sum_{i=1}^n \alpha_i y_i = 0$$

Add WeChat powcoder

Then use the obtained α_i 's to solve for the weights and bias

$$\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \qquad b = y_i - \mathbf{w}^T \mathbf{x}_i \quad \forall i$$

<https://powcoder.com>
Prediction on Test Example
[Assignment Project Exam Help](#)

Now we have the solution for the weights and bias

$$\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \quad b = y_i - \mathbf{w}^T \mathbf{x}_i \quad \forall i$$

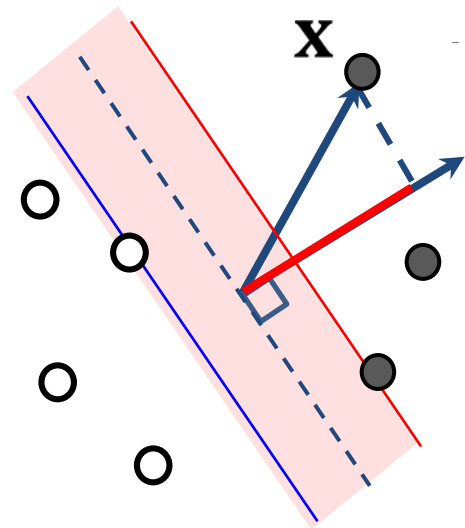
[Assignment Project Exam Help](#)

Given a new input example \mathbf{x} , classify it as

[Add WeChat powcoder](#)

$$\begin{aligned} &+1 \text{ if } \mathbf{w}^T \mathbf{x} + b \geq 1, \text{ or} \\ &-1 \text{ if } \mathbf{w}^T \mathbf{x} + b \leq -1 \end{aligned}$$

In practice, predict $y = \text{sign}[\mathbf{w}^T \mathbf{x} + b]$



<https://powcoder.com>

Dual vs Primal SVM

Assignment Project Exam Help

n is the number of training points, d is dimension of \mathbf{x} , \mathbf{w}

Add WeChat powcoder

Primal problem: for $\mathbf{w} \in \mathbb{R}^d$, hyperparameter C , the unconstrained

$$\min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i)$$

<https://powcoder.com>

Dual problem: for $\alpha \in \mathbb{R}^n$

Add WeChat powcoder

$$L = \max_{\alpha_i \geq 0} \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j (\mathbf{x}_i \cdot \mathbf{x}_j) \right\} \quad \text{s.t.} \quad \alpha_i \geq 0; \quad \sum_{i=1}^n \alpha_i y_i = 0$$

- Efficiency: need to learn d parameters for primal, n for dual

<https://powcoder.com>

Dual vs Primal SVM

Assignment Project Exam Help

Add WeChat powcoder

- Dual: quadratic programming problem in which we optimize a quadratic function of \mathbf{a} subject to a set of inequality constraints

Assignment Project Exam Help

- The solution to a quadratic programming problem in d variables in general has computational complexity that is $O(d^3)$

<https://powcoder.com>

Add WeChat powcoder

- For a fixed set of basis functions whose number d is smaller than the number n of data points, the move to the dual problem appears disadvantageous.
- However, it allows the model to be reformulated using kernels which allow *infinite* feature spaces (more on this later)

<https://powcoder.com>

Dual vs Primal SVM

Assignment Project Exam Help

Add WeChat powcoder

- Most of the SVM literature and software solves the Lagrange dual problem formulation
- Why prefer solving the dual problem over the primal?
 - provides a way to deal with constraints
 - expresses solution in terms of dot products of data points, allowing kernels
 - historical reasons

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

For an in-depth discussion refer to

<http://olivier.chapelle.cc/pub/neco07.pdf> (optional reading)

<https://powcoder.com>

Support Vectors

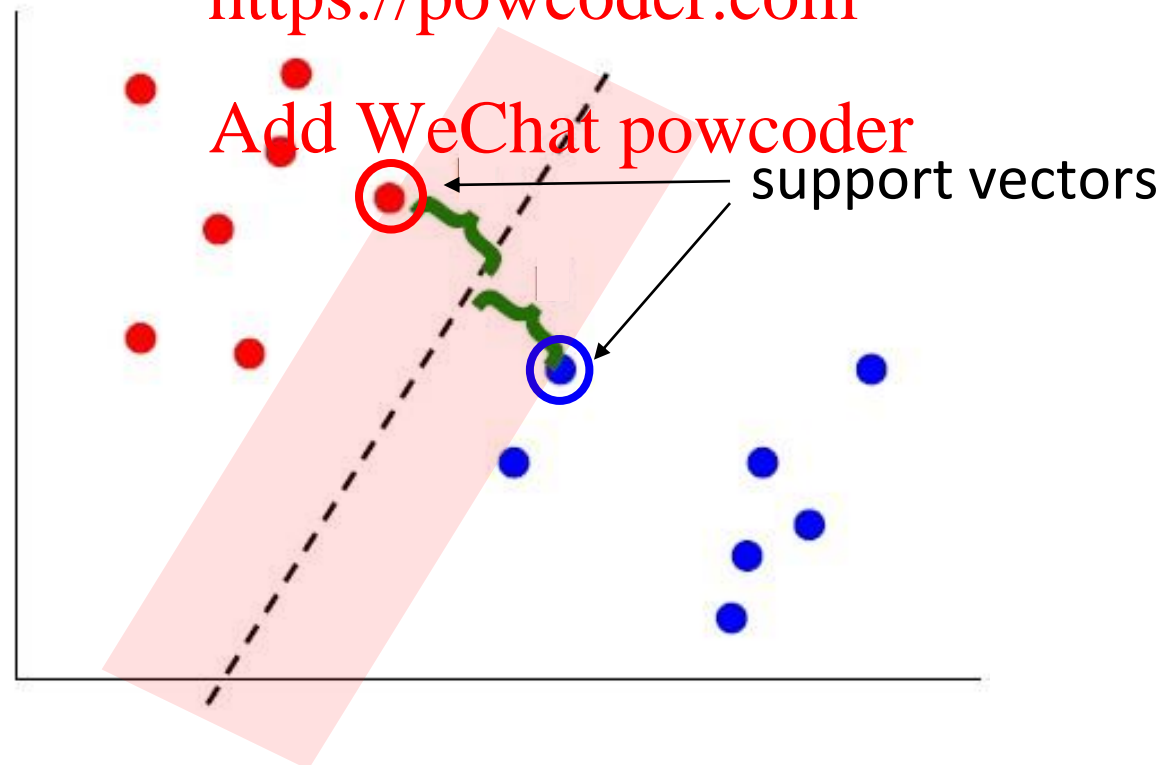
Assignment Project Exam Help

Only a small subset of α_i 's will be non-zero, and the corresponding \mathbf{x}_i 's are the **support vectors** \mathbf{S}

$$y = \text{sign}[b + \mathbf{x} \cdot (\sum_{i=1}^n y_i \alpha_i \mathbf{x}_i)] = \text{sign}[b + \mathbf{x} \cdot (\sum_{i \in \mathbf{S}} y_i \alpha_i \mathbf{x}_i)]$$

<https://powcoder.com>

Add WeChat powcoder



<https://powcoder.com> Summary of Linear SVM

Assignment Project Exam Help

- Binary and linear separable classification (regression possible too)
- Linear classifier with maximal margin

Add WeChat powcoder

- Training SVM by maximizing

$$\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

Assignment Project Exam Help

- Subject to $\alpha_i \geq 0$; $\sum_{i=1}^n \alpha_i y_i = 0$

<https://powcoder.com>

- Weights: $\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$

Add WeChat powcoder

- Only a small subset of α_i 's will be nonzero, and the corresponding \mathbf{x}_i 's are the support vectors \mathbf{S}
- Prediction on a new example:

$$y = \text{sign}[b + \mathbf{x} \cdot (\sum_{i=1}^n y_i \alpha_i \mathbf{x}_i)] = \text{sign}[b + \mathbf{x} \cdot (\sum_{i \in \mathbf{S}} y_i \alpha_i \mathbf{x}_i)]$$

<https://powcoder.com>

Next Class

Assignment Project Exam Help

Add WeChat powcoder

Support Vector Machines II

non-separable data; slack variables; kernels;
multiclass SVM

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Reading: Bishop Ch 6.1-6.2, Ch 7.1.3