

# Announcements

**Reminder:** ps5 out, due Thursday 11/5

- pset 4 grades up on blackboard by Monday 11/9

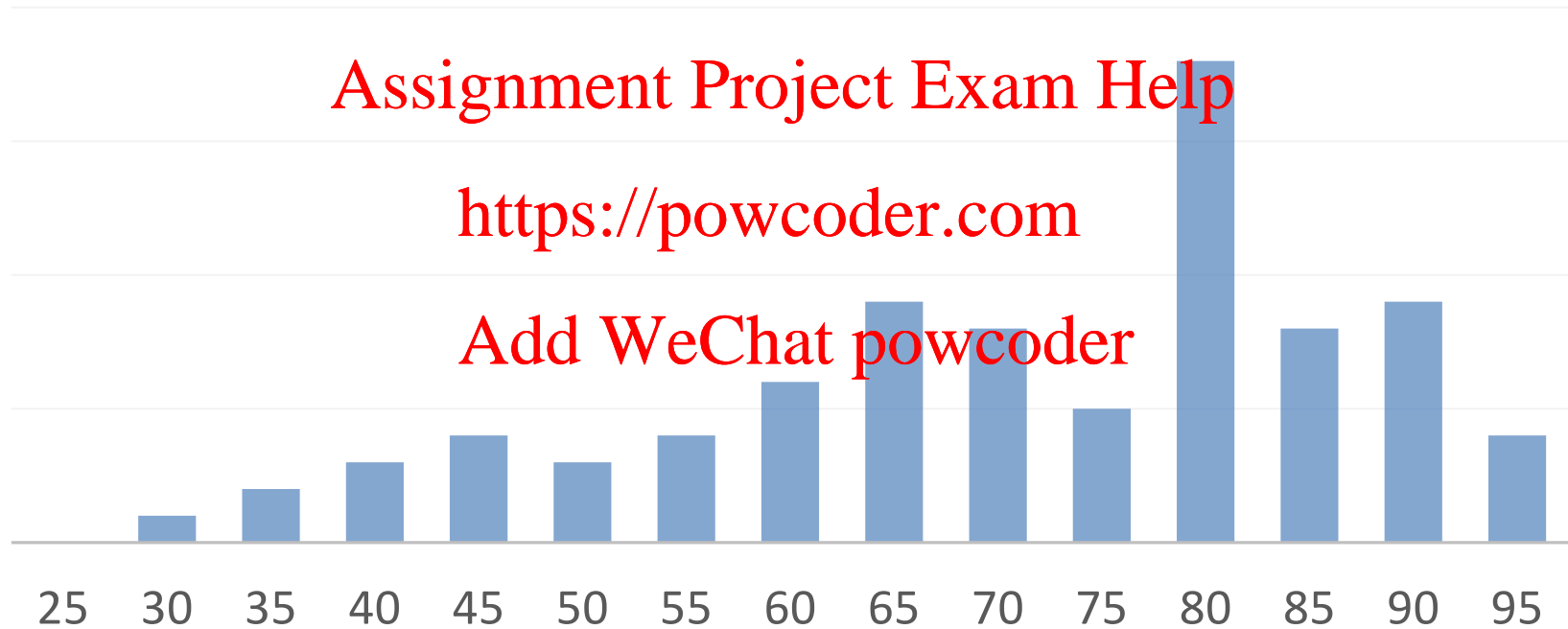
Assignment Project Exam Help

<https://powcoder.com>

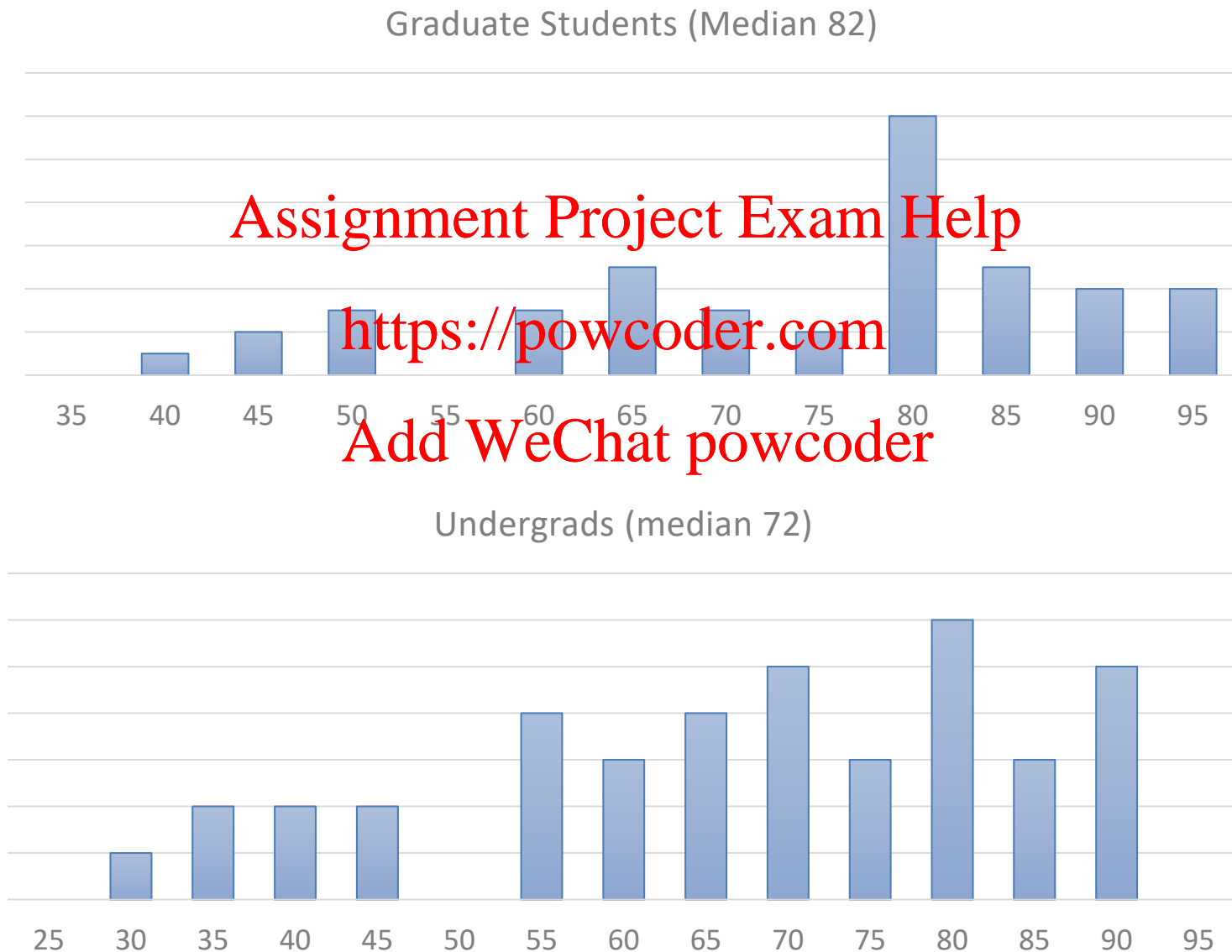
Add WeChat powcoder

# Midterm grades out!

Unweighted midterm grades (Median = 78)



# Graduate students did better overall



# Class grading

- 20% midterm
- 20% final
- 15% class challenge
- 45% homeworks

Assignment Project Exam Help

<https://powcoder.com>

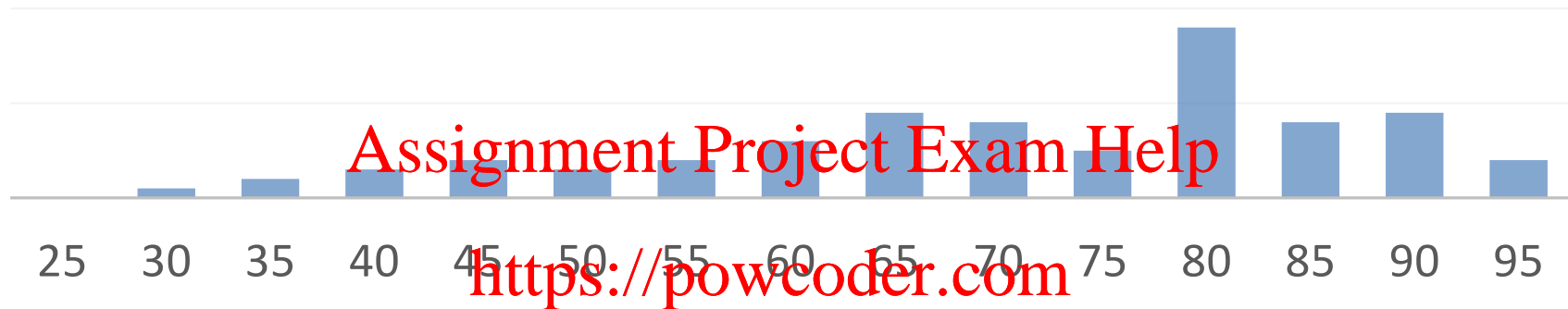
Add WeChat powcoder

Assume student gets 72% on midterm and final,  
85% on homework/challenge= ~80% (B-)

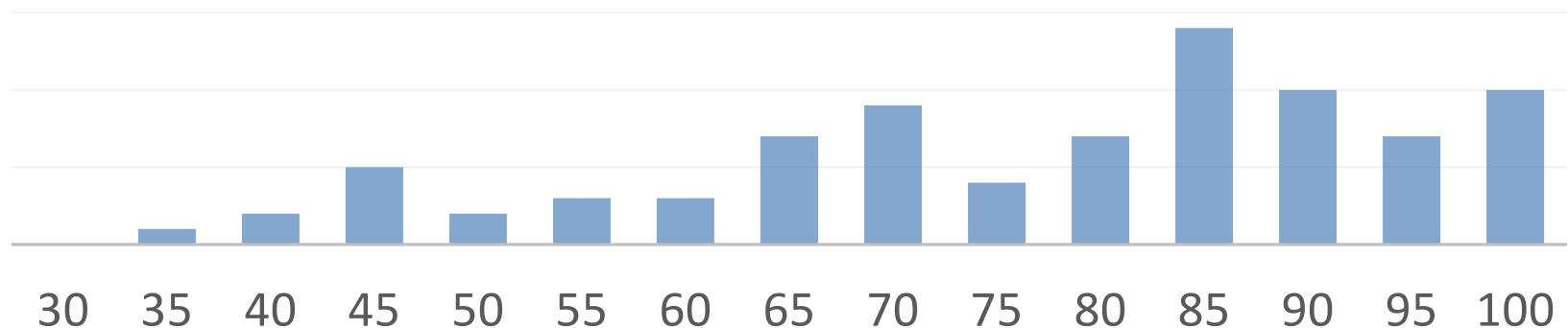
72% on midterm, 85% final, 95%  
homework/challenge=~88% (B+)

# Two questions < 50% points awarded, retroactively make them bonus points

Unweighted midterm grades (Median = 78)



Weighted Midterm (Median = 85)



New median for graduates: 89%  
New median for undergraduates: 78%

# Class grading

- 20% midterm
- 20% final
- 15% class challenge
- 45% homeworks

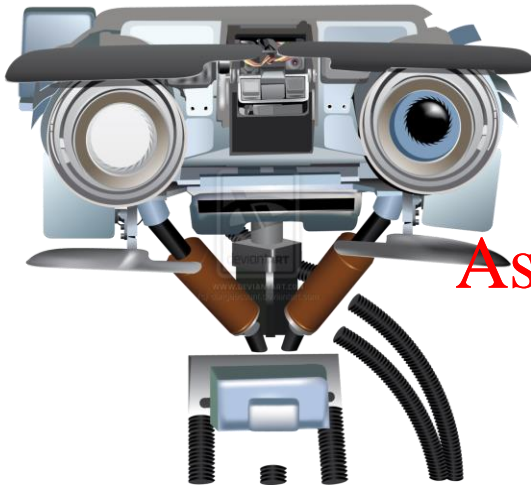
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Assume student gets 78% on midterm and final,  
85% on homework/challenge= ~82% (B-/B)

78% on midterm, 88% final, 95%  
homework/challenge=~90% (A-)



Assignment Project Exam Help

Soft-margin SVM  
<https://powcoder.com>

---

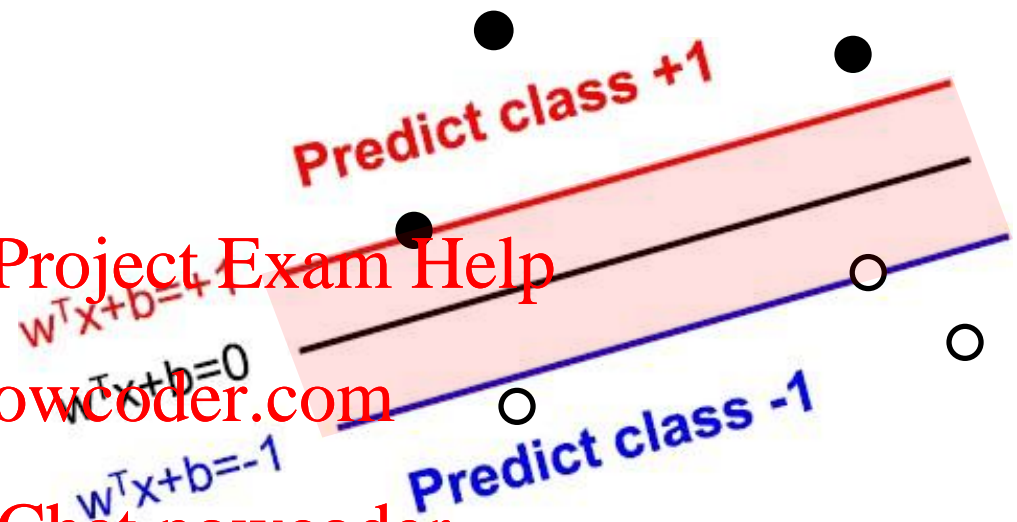
Add WeChat powcoder

# Max Margin Classifier

“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



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Classify as +1

if

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

Classify as -1

if

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

Undefined

if

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



# Dual vs Primal SVM

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

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

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

Assignment Project Exam Help  
<https://powcoder.com>

**Dual problem:** for  $\alpha \in \mathbb{R}^n$   
 Add WeChat powcoder

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

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

# What if data is not linearly separable?

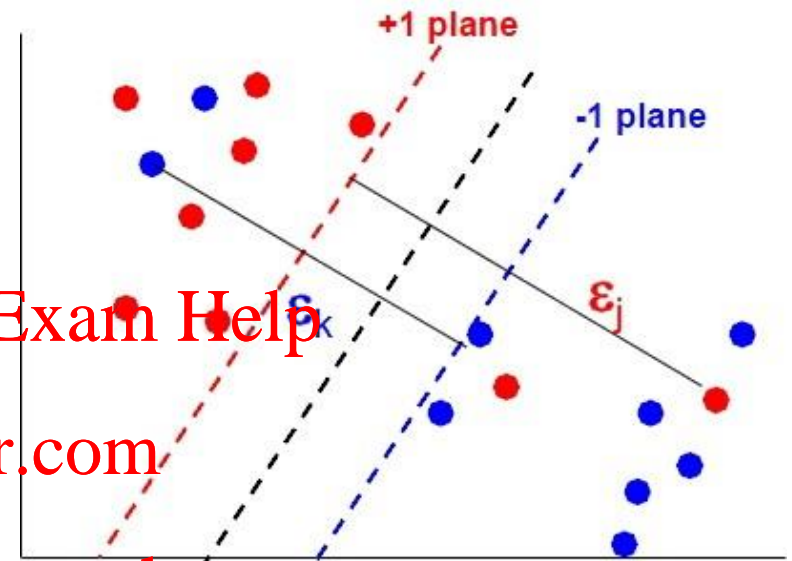
- Introduce **slack variables**  $\xi_i$

$$\min \left[ \frac{1}{2} \|\mathbf{w}\|^2 + \lambda \sum_{i=1}^n \xi_i \right]$$

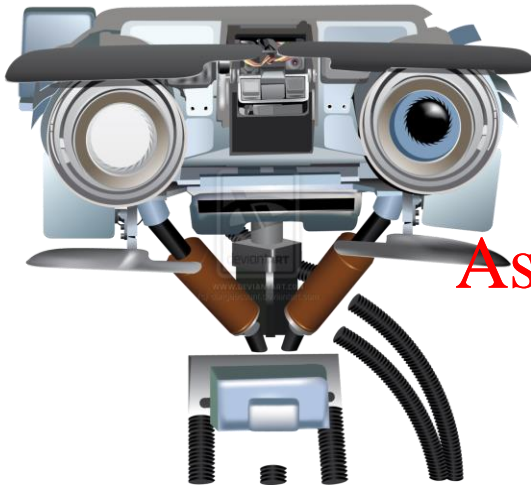
subject to constraints (for all  $i$ ):

$$y_i (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$$

$$\xi_i \geq 0$$



- Example lies on wrong side of hyperplane:  $\xi_i > 1 \Rightarrow \sum_i \xi_i$  is upper bound on number of training errors
- $\lambda$  trades off training error versus model complexity
- This is known as the **soft-margin** extension



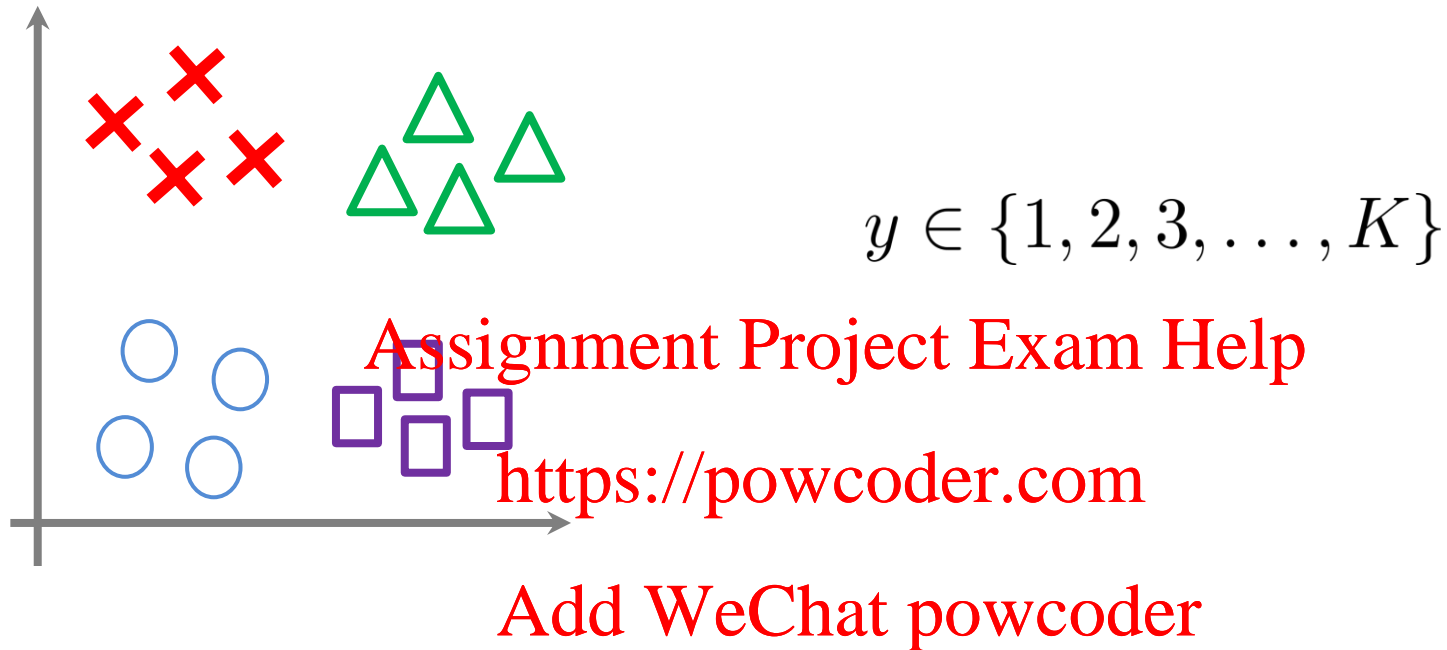
Assignment Project Exam Help

Multi-class SVMs  
<https://powcoder.com>

---

Add WeChat powcoder

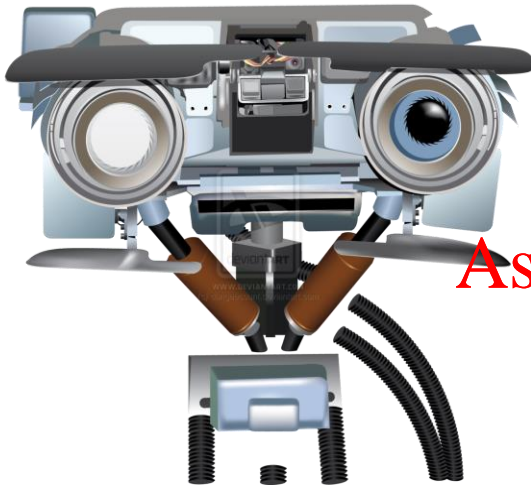
# Multi-class classification



Many SVM packages already have built-in multi-class classification functionality.

Otherwise, use one-vs.-all method. (Train  $K$  SVMs, one to distinguish class  $i$  from the rest), for  $i = 1, \dots, K$ , get  $\mathbf{w}^{(1)}, b^{(1)}, \dots, \mathbf{w}^{(K)}, b^{(K)}$

Pick class  $y = i$  with largest score  $\mathbf{w}^{(i)T} \mathbf{x} + b^{(i)}$



Assignment Project Exam Help  
**Kernel SVM**  
<https://powcoder.com>

---

Add WeChat powcoder

# Non-linear decision boundaries

- Note that both the learning objective and the decision function depend only on dot products between patterns

$$L = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j (\mathbf{x}_i \cdot \mathbf{x}_j) \quad y = \text{sign}[b + \mathbf{x} \cdot (\sum_{i=1}^n y_i \alpha_i \mathbf{x}_i)]$$

- How to form non-linear decision boundaries in input space?

- Basic idea: <https://powcoder.com>

- Map data into feature space  $\mathbf{x} \rightarrow \phi(\mathbf{x})$
- Replace dot products between inputs with feature points  
 $\mathbf{x}_i \cdot \mathbf{x}_j \rightarrow \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$
- Find linear decision boundary in feature space

- Problem: what is a good feature function  $\phi(\mathbf{x})$ ?

# Kernel Trick

- **Kernel trick**: dot-products in feature space can be computed as a kernel function

$$\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j) = K(\mathbf{x}_i, \mathbf{x}_j)$$

- Idea: work directly on  $\mathbf{x}$ , avoid having to compute  $\phi(\mathbf{x})$

<https://powcoder.com>

- Example:

**Add WeChat powcoder**

$$\begin{aligned} K(\mathbf{a}, \mathbf{b}) &= (\mathbf{a} \cdot \mathbf{b})^3 = ((a_1, a_2) \cdot (b_1, b_2))^3 \\ &= (a_1 b_1 + a_2 b_2)^3 \\ &= a_1^3 b_1^3 + 3a_1^2 b_1^2 a_2 b_2 + 3a_1 b_1 a_2^2 b_2^2 + a_2^3 b_2^3 \\ &= (a_1^3, \sqrt{3}a_1^2 a_2, \sqrt{3}a_1 a_2^2, a_2^3) \cdot (b_1^3, \sqrt{3}b_1^2 b_2, \sqrt{3}b_1 b_2^2, b_2^3) \\ &= \phi(\mathbf{a}) \cdot \phi(\mathbf{b}) \end{aligned}$$

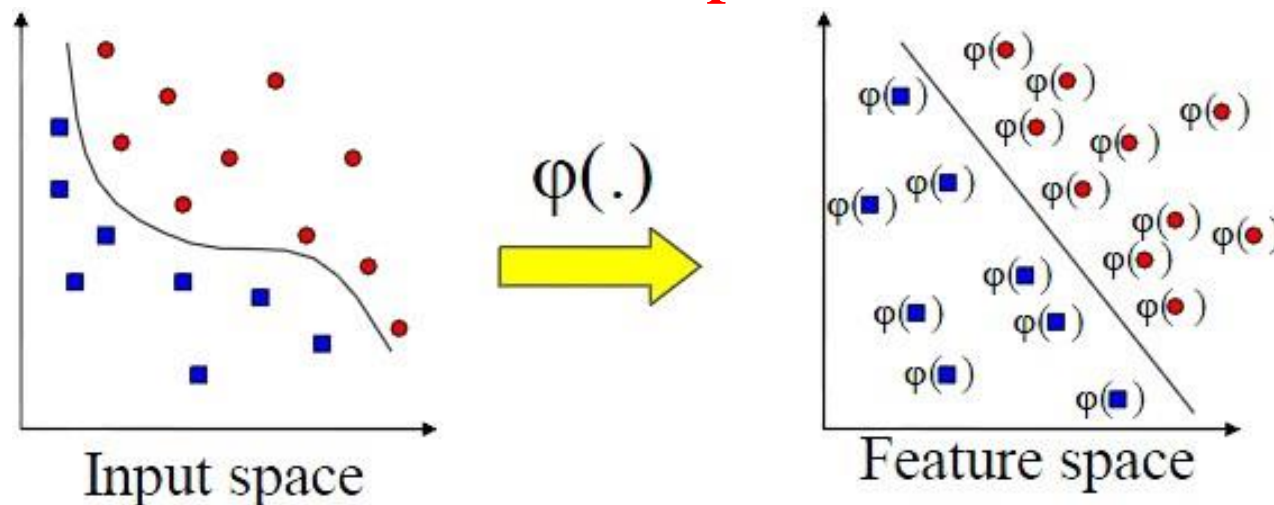
# Input transformation

Mapping to a feature space can produce problems:

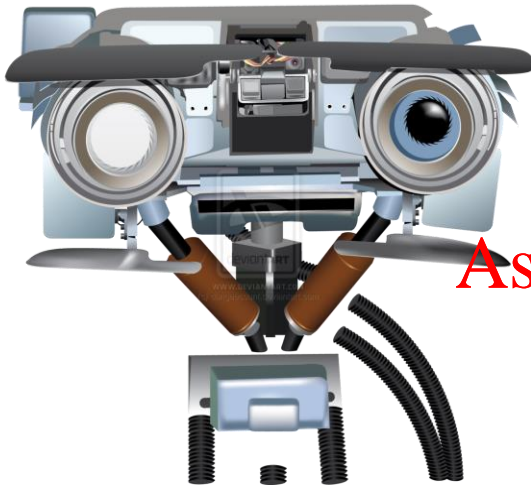
- High computational burden due to high dimensionality
- Many more parameters

SVM solves these two issues simultaneously

- Kernel trick produces efficient classification
- Dual formulation only assigns parameters to samples, not features







Assignment Project Exam Help

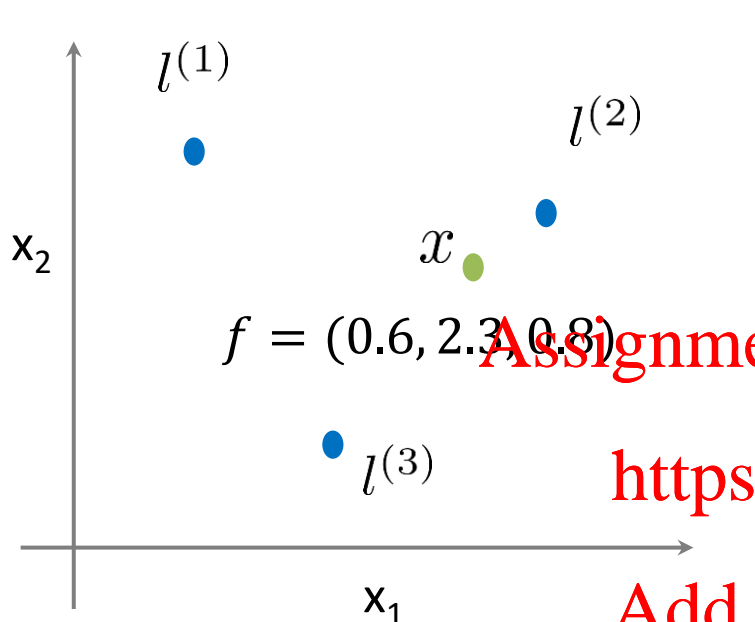
Kernels

<https://powcoder.com>

---

Add WeChat powcoder

# Kernels as Similarity Functions



Given  $x$ , compute new feature depending on proximity to “landmarks”  $l^{(1)}, l^{(2)}, l^{(3)}$

Example: Gaussian (RBF) kernel

$$f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

If  $x \approx l^{(1)}$  :

similarity is high

If  $x$  is far from  $l^{(1)}$  :

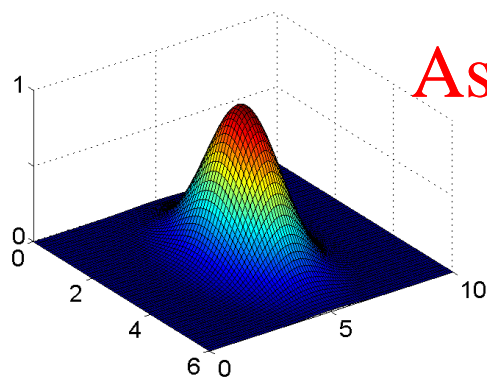
similarity is low

Predict label “1” when  $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$

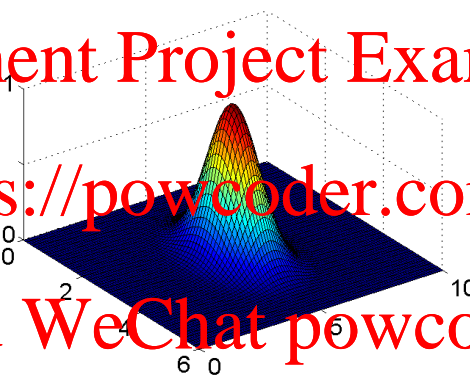
**Example:**

$$l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}, \quad f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

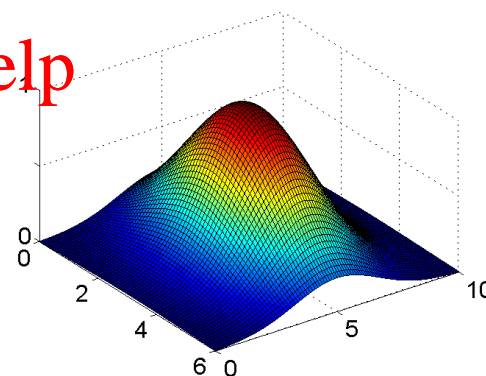
$$\sigma^2 = 1$$



$$\sigma^2 = 0.5$$



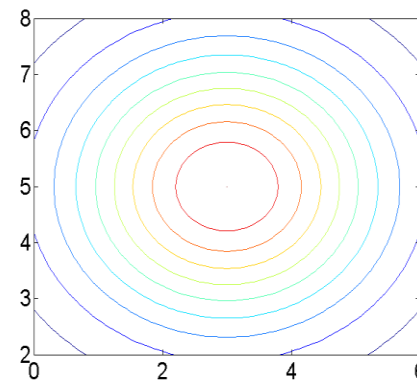
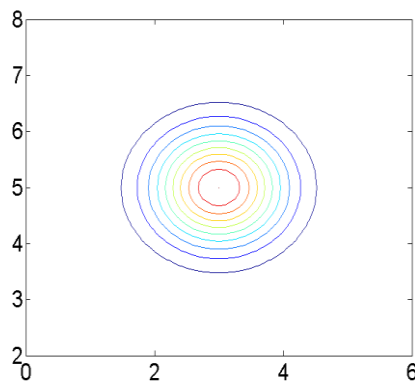
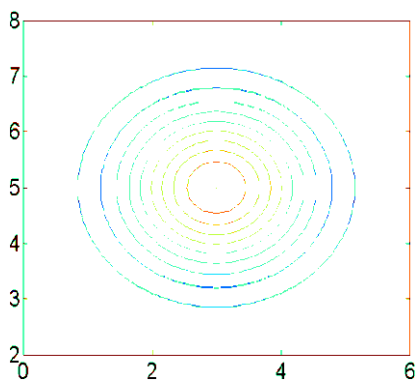
$$\sigma^2 = 3$$



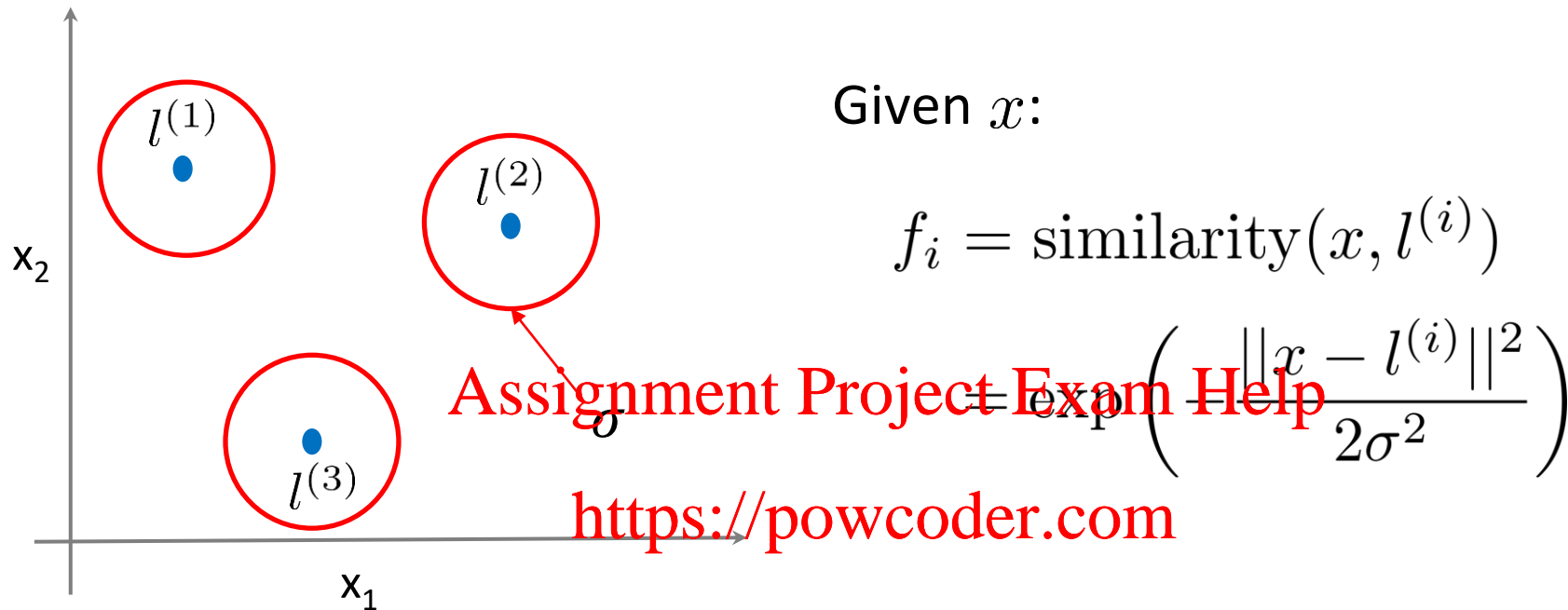
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



## Landmarks for Gaussian kernel



**Add WeChat powcoder**

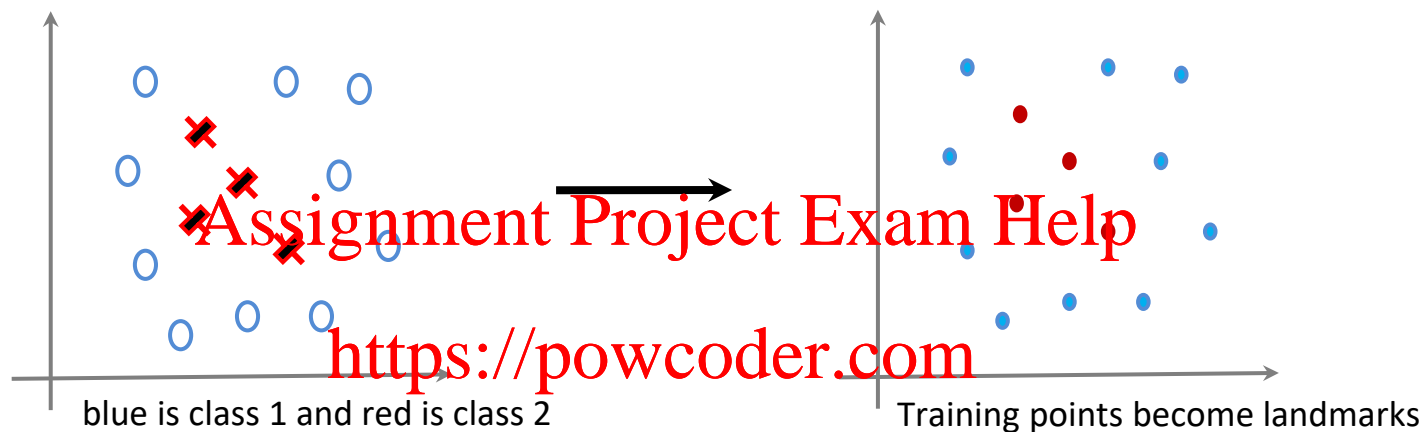
Predict  $y = 1$  if  $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$

So the new features  $f$  measure how close the example is to each “landmark” point

Where do the landmarks come from?

## Landmarks for Gaussian kernel

Where do  $l^{(1)}, l^{(2)}, l^{(3)}, \dots$  come from? They are the training points!



Add WeChat powcoder

So the "landmarks" are points we can use to compute a new feature representation for a point  $x$ , by representing it as the similarity to each landmark point (measured using a Gaussian centered at the landmark)

In SVMs with RBF (Gaussian) kernels, we place a Gaussian centered at **each** training point to compute the nonlinear features. This is equivalent to using all of the training data as landmarks.

## SVM with Kernels

Given  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$ ,  
choose  $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(m)} = x^{(m)}$ .

Given example  $x$ :

$$f_1 = \text{similarity}(x, l^{(1)})$$

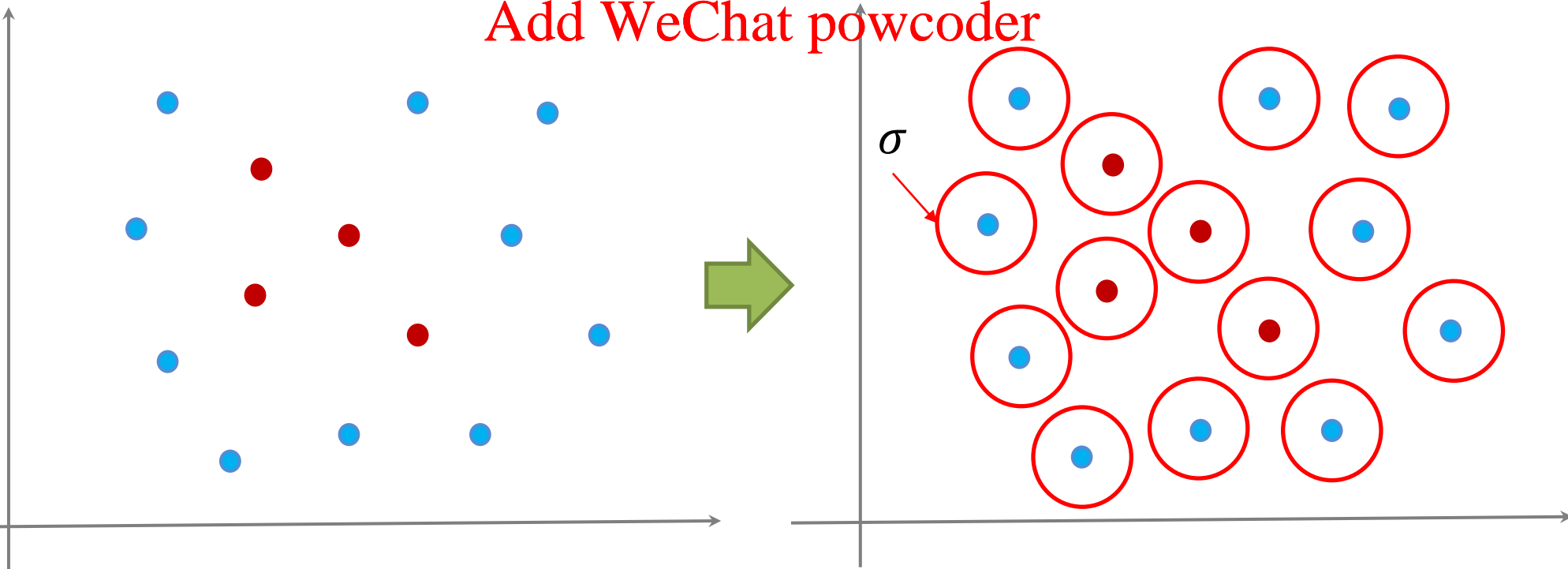
$$f_2 = \text{similarity}(x, l^{(2)})$$

...

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



# Kernels

Examples of kernels (kernels measure similarity):

1. Polynomial  $K(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1 \cdot \mathbf{x}_2 + 1)^2$
2. Gaussian  $K(\mathbf{x}_1, \mathbf{x}_2) = \exp(-\|\mathbf{x}_1 - \mathbf{x}_2\|^2 / 2\sigma^2)$
3. Sigmoid  $K(\mathbf{x}_1, \mathbf{x}_2) = \tanh(\kappa(\mathbf{x}_1 \cdot \mathbf{x}_2) + a)$

Assignment Project Exam Help

<https://powcoder.com>

Each kernel computation corresponds to dot product calculation for particular mapping  $\phi(\mathbf{x})$ . Implicitly maps to high dimensional space

Add WeChat powcoder

Why is this useful?

1. Rewrite training examples using more complex features
2. Dataset not linearly separable in original space may be linearly separable in higher dimensional space

# Classification with non-linear SVMs

Non-linear SVM using kernel function  $K()$ :

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

Maximize  $L_K$  w.r.t.  $\{\alpha\}$ , under constraints  $\alpha \geq 0$

<https://powcoder.com>

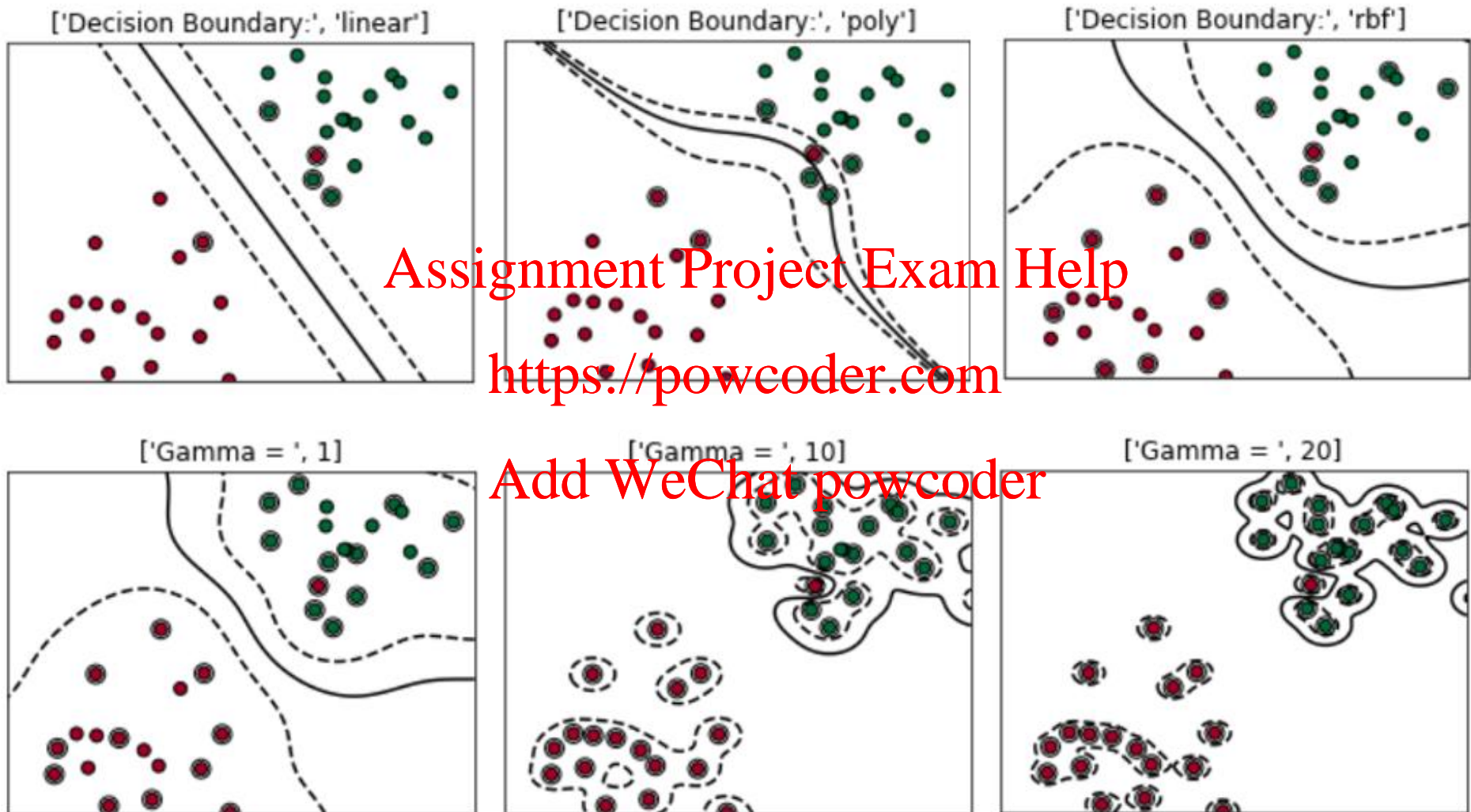
Unlike linear SVM, cannot express  $w$  as linear combination of support vectors – now must retain the support vectors to classify new examples

Final decision function:

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



# Decision boundary in Gaussian kernel SVM



$$\text{Gamma} = 1/\sigma^2$$

# Kernel SVM Summary

## Advantages:

- Kernels allow very flexible hypotheses
- Poly-time exact optimization methods rather than approximate methods
- Soft-margin extension permits mis-classified examples
- Variable-sized hypothesis space
- Excellent results (1.1% error rate on handwritten digits vs. LeNet's 0.9%)

## Disadvantages:

- Must choose kernel hyperparameters
- Very large problems computationally intractable
- Batch algorithm

# Kernel Functions

Mercer's Theorem (1909): any reasonable kernel corresponds to some feature space

Reasonable means that the Gram matrix is positive definite

$$K_p = K(\mathbf{x}_p, \mathbf{x}_q)$$

Feature space can be very large, e.g., polynomial kernel  $(1 + \mathbf{x}_i + \mathbf{x}_j)^d$  corresponds to feature space exponential in  $d$

Linear separators in these super high-dim spaces correspond to highly nonlinear decision boundaries in input space

# Kernelizing

A popular way to make an algorithm more powerful is to develop a kernelized version of it

- We can rewrite a lot of algorithms to be defined only in terms of inner product <https://powcoder.com>
- For example: k-nearest neighbors

$$\mathbf{z} = \varphi(\mathbf{x})$$

$$(\mathbf{z}_i - \mathbf{z}_j)^2 = K(\mathbf{x}_i, \mathbf{x}_i) + K(\mathbf{x}_j, \mathbf{x}_j) - 2K(\mathbf{x}_i, \mathbf{x}_j)$$

# Techniques for constructing valid kernels

Given valid kernels  $k_1(\mathbf{x}, \mathbf{x}')$  and  $k_2(\mathbf{x}, \mathbf{x}')$ , the following new kernels will also be valid:

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \quad (6.13)$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \quad (6.14)$$

$$k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.15)$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.16)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}') \quad (6.17)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}') \quad (6.18)$$

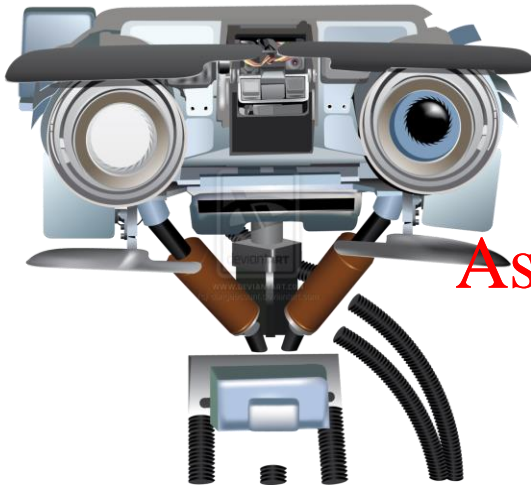
$$k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}')) \quad (6.19)$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{A} \mathbf{x}' \quad (6.20)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.21)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.22)$$

where  $c > 0$  is a constant,  $f(\cdot)$  is any function,  $q(\cdot)$  is a polynomial with nonnegative coefficients,  $\phi(\mathbf{x})$  is a function from  $\mathbf{x}$  to  $\mathbb{R}^M$ ,  $k_3(\cdot, \cdot)$  is a valid kernel in  $\mathbb{R}^M$ ,  $\mathbf{A}$  is a symmetric positive semidefinite matrix,  $\mathbf{x}_a$  and  $\mathbf{x}_b$  are variables (not necessarily disjoint) with  $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$ , and  $k_a$  and  $k_b$  are valid kernel functions over their respective spaces.



Assignment Project Exam Help

Summary of SVMs  
<https://powcoder.com>

---

Add WeChat powcoder

# Summary

## Software:

- A list of SVM implementations can be found at <http://www.kernel-machines.org/software.html>
- Some implementations (such as LIBSVM) can handle multi-class classification
- SVMLight is among the earliest implementations
- Several Matlab toolboxes for SVM are also available

## Key points:

- Difference between logistic regression and SVMs
- Maximum margin principle
- Target function for SVMs
- Slack variables for mis-classified points
- Kernel trick allows non-linear generalizations

# History of SVMs

- The original SVM algorithm was invented by [Vladimir Vapnik](#) and [Alexey Chervonenkis](#) in 1963.
- In 1992, Bernhard E. Boser, Isabelle M. Guyon and [Vladimir Vapnik](#) suggested a way to create nonlinear classifiers by applying the [kernel trick](#) to maximum-margin hyperplanes. [\[13\]](#)
- The soft margin was proposed by [Corinna Cortes](#) and Vapnik in 1993 and published in 1995. [\[1\]](#)
- SVMs were very popular in the 90's-00's until neural networks took over circa 2012



# Next Class

## Reinforcement Learning I

reinforcement learning; Markov Decision  
Process (MDP); policies, value functions, Q-  
learning

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder