Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

THE UNIVERSITY OF
MELBOURNE

- Assuming the data is linearly separable

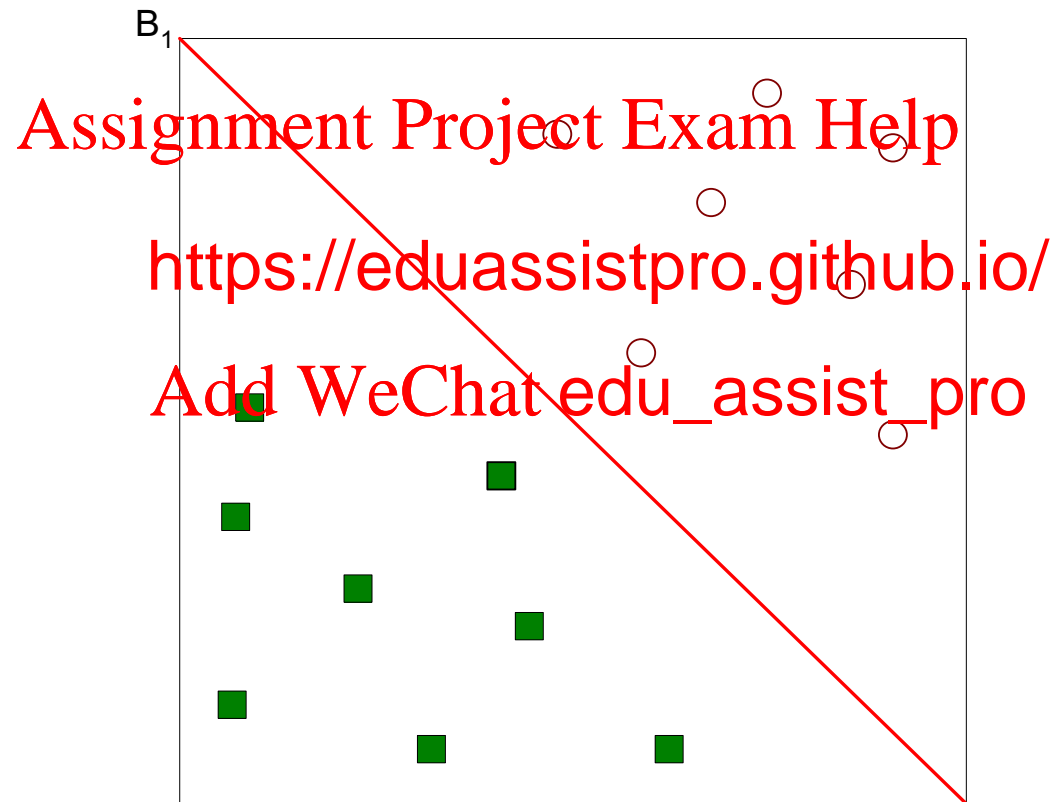- Aim: find a linear hyperplane (decision boundary) that will separate the data
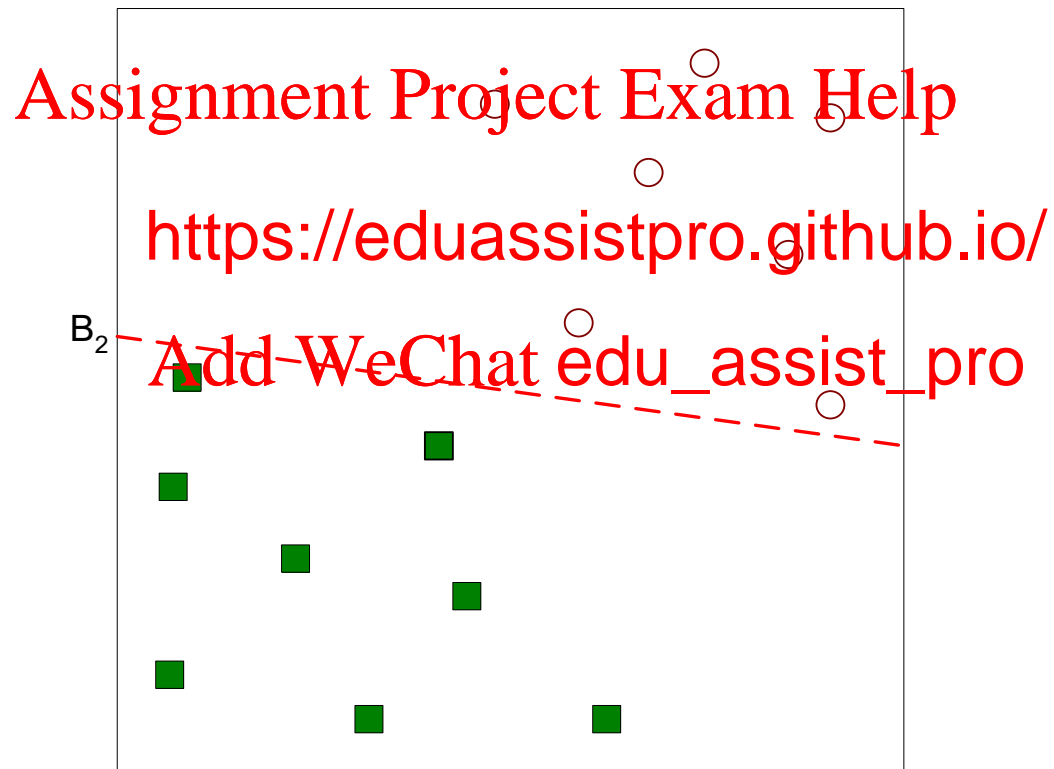
Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

THE UNIVERSITY OF MELBOURNE

- One Possible Solution

B$_1$

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

- Another Possible Solution

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

$B_2$

THE UNIVERSITY OF MELBOURNE

- Other Possible Solutions



Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

$B_2$

- Which one is better? B1 or B2?
- How do you define better?

- Find hyperplane maximises the margin => B1 is better than B2
- Margin: sum of shortest distances from the planes to the positive/negative samples

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

- Small margin separating planes:
  - are more fragile to noise
  - may over-fit the data

- Large margin separating planes:
  - are more robust to n
  - From statistical lear                        nes
    generalises better to unseen data

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

$\{\mathbf{x}_i, y_i\}$ where $i = 1 \ldots L, y_i \in \{-1, 1\}, \mathbf{x}_i \in \mathbb{R}^D$

This hyperplane can be described by $\mathbf{x} \cdot \mathbf{w} + b = 0$ where:

- $\mathbf{w}$ is normal to the hyperplane.

- $\frac{b}{\|\mathbf{w}\|}$ is the perpendicular distance from the hyperplane to the origin.

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

$\{\mathbf{x}_i, y_i\}$ where $i = 1 \ldots L, y_i \in \{-1, 1\}, \mathbf{x}_i \in \mathbb{R}^D$
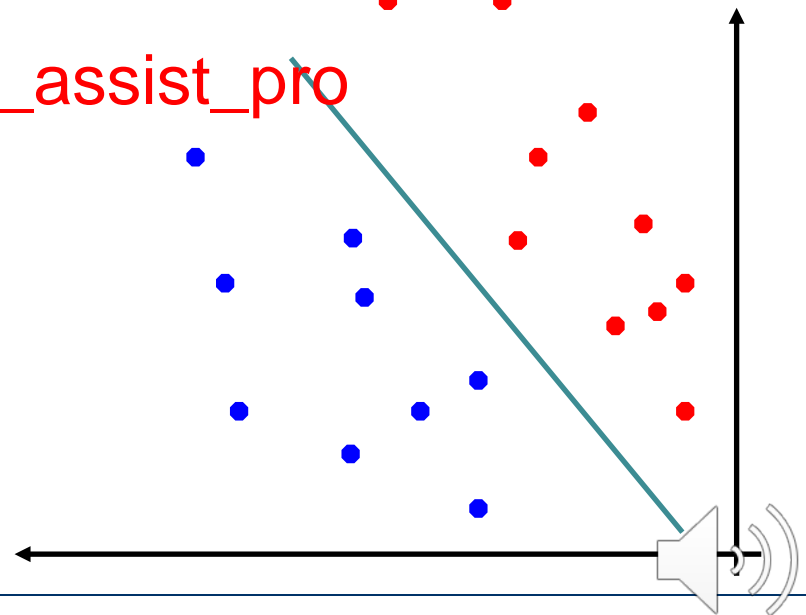
This hyperplane can be described by $\mathbf{x} \cdot \mathbf{w} + b = 0$ where:

- $\mathbf{w}$ is normal to the hyperplane.

- $\frac{b}{\|\mathbf{w}\|}$ is the perpendicular distance from the hyperplane to the origin.

Assignment Project Exam Help

https://eduassistpro.github.io/

$\mathbf{x} \cdot \mathbf{w} + b \geq 0$

**Classification r**

Add WeChat edu_assist_pro

$$f(\mathbf{x}) = \text{sign}(\mathbf{x} \cdot \mathbf{w} + b) = \begin{cases} +1 & \text{if } \mathbf{x} \cdot \mathbf{w} + b \geq 0 \\ -1 & \text{if } \mathbf{x} \cdot \mathbf{w} + b < 0 \end{cases}$$

Find $\mathbf{w}$ and $b$ such that:

$$\mathbf{x}_i \cdot \mathbf{w} + b \geq 0 \text{ for } y_i = +1$$
$$\mathbf{x}_i \cdot \mathbf{w} + b < 0 \text{ for } y_i = -1$$
$$\text{for all } i = 1 \ldots L$$

$\mathbf{x} \cdot \mathbf{w} + b < 0$

**Training objective**

$B_1$

$\mathbf{x} \cdot \mathbf{w} + b = 0$

$\mathbf{x} \cdot \mathbf{w} + b = -1$

**Support Vectors**

Assignment Project Exam Help

https://eduassistpro.github.io/

$\mathbf{x} \cdot \mathbf{w} + b = +1$

Add WeChat edu_assist_pro

11

$b_{12}$    $\text{Margin} = \dfrac{2}{\|\mathbf{w}\|}$

**Requirement for margin:**

$$\mathbf{x}_i \cdot \mathbf{w} + b \geq +1 \quad \text{for } y_i = +1$$
$$\mathbf{x}_i \cdot \mathbf{w} + b \leq -1 \quad \text{for } y_i = -1$$

**Margin**

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

**Note that:**

$$\mathbf{x}_i \cdot \mathbf{w} + b \geq +1 \quad \text{for } y_i = +1$$
$$\mathbf{x}_i \cdot \mathbf{w} + b \leq -1 \quad \text{for } y_i = -1$$

These equations can be combined into:

$$y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \geq 0 \ \forall_i$$

(1)   $$\max \frac{2}{\|\mathbf{w}\|} \quad \text{s.t.} \quad y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \geq 0 \quad \forall i$$

(2)

Assignment Project Exam Help

https://eduassistpro.github.io/

(3)

Add WeChat edu_assist_pro

$$\min \frac{1}{2}\|\mathbf{w}\|^2 \quad \text{s.t.} \quad y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \geq 0 \quad \forall_i$$

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

- For linearly separable data: a max-margin solution is **guaranteed** to exist
- For non- linearly separable data: a solution does not exist

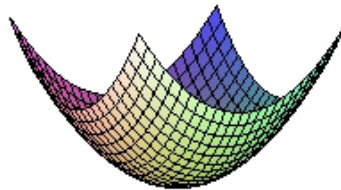$$\min \frac{1}{2}\|\mathbf{w}\|^2 \quad \text{s.t.} \quad y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \geq 0 \quad \forall_i$$

- Need to optimize a *quadratic* function subject to *linear* constraints.
- Convex quadratic optimization problem
- Convex objective: any local minimum

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

**Convex**　　　**Non Convex**

**Primal problem**: solve for **w** and b

$$\min \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \geq 0 \ \ \forall_i$$

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

**Primal problem**: solve for **w** and b

$$\min \frac{1}{2} \|\mathbf{w}\|^2 \qquad \text{s.t.} \qquad y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \geq 0 \;\; \forall_i$$

*Equivalent **dual problem*** formulation: solve for $a_1 \ldots a_L$ Lagrange multipliers for each data point

*More convenient to solve*

See Ref. [1] for derivation

- Given a solution $\alpha_1 \dots \alpha_L$ to the dual problem, solution to the primal is:

$$\mathbf{w} = \Sigma \alpha_i y_i \mathbf{x}_i \qquad b = y_k - \Sigma \alpha_i y_i \mathbf{x}_i{}^\mathsf{T} \mathbf{x}_k \quad \text{for any } \alpha_k > 0$$

- Each non-zero $\alpha_i$ indicates that corresponding $\mathbf{x}_i$ is a support vector.
- Then the classifying fun                                    ed $\mathbf{w}$ explicitly):

$$f(\mathbf{x}) = \Sigma \alpha_i y_i \mathbf{x}_i{}^\mathsf{T} \mathbf{x} + b$$

- Notice that it relies on an *inner product* between the test point $\mathbf{x}$ and the support vectors $\mathbf{x}_i$ – we will return to this later.
- Also keep in mind that solving the optimization problem involved computing the inner products $\mathbf{x}_i{}^\mathsf{T}\mathbf{x}_j$ between all training points.

THE UNIVERSITY OF
MELBOURNE

- Classification function:

$$f(\mathbf{x}) = \Sigma \alpha_i y_i \mathbf{x}_i^\mathbf{T} \mathbf{x} + b$$

Linear
SVM

- Only support vectors matter, other training examples are ignorable.

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

- What if the training set is mostly, but not exactly, linearly separable?



Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

*The (hard) linear SVM problem is **infeasible** here.*

- **Slack variables** $\xi_i$ can be added to allow misclassification of difficult or noisy examples, resulting margin called *soft*.



Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

- The old formulation (hard SVM):

Find $\mathbf{w}$ and b such that
$\Phi(\mathbf{w}) = \mathbf{w}^T\mathbf{w}$ is minimized
and for all $(\mathbf{x}_i, y_i)$, $i=1..L$ : $y_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1$

Assignment Project Exam Help

- Modified formulation inc https://eduassistpro.github.io/ oft SVM):

Find $\mathbf{w}$ and b such that Add WeChat edu_assist_pro
$\Phi(\mathbf{w}) = \mathbf{w}^T\mathbf{w} + C\Sigma\xi_i$ is minimized
and for all $(\mathbf{x}_i, y_i)$, $i=1..L$ : $y_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1 - \xi_{i,}$ , $\xi_i \geq 0$

- **Parameter C** can be viewed as a way to control overfitting: it "trades off" the relative importance of maximizing the margin and fitting the training data.

- Dual problem is identical to separable case:

Find $\alpha_1...\alpha_L$ such that

$\mathbf{Q}(\boldsymbol{\alpha}) = \Sigma\alpha_i - \frac{1}{2}\Sigma\Sigma\alpha_i\alpha_j y_i y_j \mathbf{x}_i^T\mathbf{x}_j$ is maximized and

(1) $\Sigma\alpha_i y_i = 0$

(2) $0 \leq \alpha_i \leq$

https://eduassistpro.github.io/

- Again, $\mathbf{x}_i$ with non-zero $\alpha_i$ will be support vector

- Solution to the primal problem is: Add WeChat edu_assist_pro

$\mathbf{w} = \Sigma\alpha_i y_i \mathbf{x}_i$

$b = y_k(1- \xi_k) - \Sigma\alpha_i y_i \mathbf{x}_i^T\mathbf{x}_k$   for any $k$ s.t. $\alpha_k > 0$

- Again, we don't need to compute $\mathbf{w}$ explicitly for classification:

$$f(\mathbf{x}) = \Sigma\alpha_i y_i \mathbf{x}_i^T\mathbf{x} + b$$

- The classifier is a *separating hyperplane*

- Most "important" training points are support vectors; they define the hyperplane.

- Quadratic optimization algorithms can identify which training points $\mathbf{x}_i$ are support vectors with non~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~ $\alpha_i$.

- Model complexity depends on #support v

- Both in the dual formulation of the problem and in the solution, training points appear only inside inner products:

Find $\alpha_1 \dots \alpha_L$ such that
$\mathbf{Q}(\boldsymbol{\alpha}) = \Sigma\alpha_i - \tfrac{1}{2}\Sigma\Sigma\alpha_i\alpha_j y_i y_j \boxed{\mathbf{x}_i^{\mathbf{T}}\mathbf{x}_j}$ is maximized and
(1) $\Sigma\alpha_i y_i = 0$
(2) $0 \leq \alpha_i \leq C$ for all $\alpha_i$

$f(\mathbf{x}) = \Sigma\alpha_i y_i \boxed{\mathbf{x}_i^{\mathbf{T}}\mathbf{x}} + b$

# Overfitting - Underfitting

**Underfitting**: model not expressive enough to capture patterns in the data

**Overfitting**: model too complicated; capture noise the data

**St right**: model res essential ns in the data

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

Linear model underfitting: model not expressive enough to capture patterns in the data

Assignment Project Exam Help

https://eduassistpro.github.io/
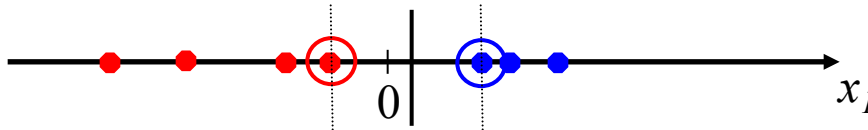
t-Margin (linear) SVM

Add WeChat edu_assist_pro

can for a small

ber of training errors

But is still a linear model

- Datasets that are linearly separable with some noise work out great:



- But what are we going to do if the dataset is just too hard?

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

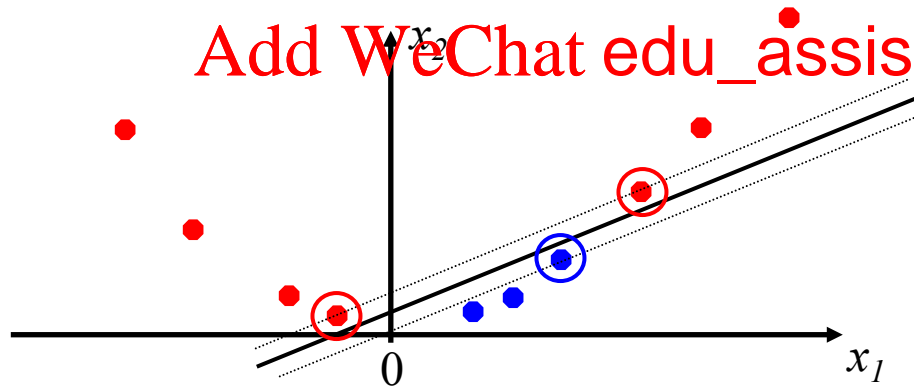- How about… mapping data to a higher-dimensional space:

# Non-linear SVMs Overview

- Turn linear SVM into a non-linear model
- By mapping the original data into a high dimensional space where the data is hopefully linearly separable

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

- General idea: the original feature space can be mapped to some higher-dimensional feature space where the training set is separable:

$$\Phi: x \rightarrow \phi(x)$$

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

- Higher-dimensional space still has *intrinsic* dimensionality *d*, but linear separators in it correspond to *non-linear* separators in original space.

Find $\alpha_1 \ldots \alpha_L$ such that
$Q(\alpha) = \Sigma \alpha_i - \tfrac{1}{2} \Sigma\Sigma \alpha_i \alpha_j y_i y_j \boxed{\mathbf{x}_i^T\mathbf{x}_j}$ is maximized and
(1) $\Sigma \alpha_i y_i = 0$
(2) $0 \leq \alpha_i \leq C$ for all $\alpha_i$

$$f(\mathbf{x}) = \Sigma \alpha_i y_i \boxed{\mathbf{x}_i^T\mathbf{x}} + b$$

Assignment Project Exam Help

- The linear SVM classifi  https://eduassistpro.github.io/  ween vectors $\mathbf{x}_i^T\mathbf{x}_j$

*(pair-wise dot products b*

Add WeChat edu_assist_pro

- If every data point is mapped into high-di                ace via some transformation $\Phi : \mathbf{x} \rightarrow \varphi(\mathbf{x})$, the inner product becomes:

$$\varphi(\mathbf{x}_i)^T\varphi(\mathbf{x}_j)$$

- Explicit mapping & Plug

$$\boldsymbol{\varphi(x_i)}^T\boldsymbol{\varphi(x_j)}$$

In place of

$$\boldsymbol{x_i}^T\boldsymbol{x_j}$$

Find $\alpha_1...\alpha_L$ such that

$\mathbf{Q(\alpha)} = \Sigma\alpha_i - \frac{1}{2}\Sigma\sigma\alpha_i\alpha_j y_i y_j \boxed{\varphi(\boldsymbol{x_i})^T \varphi(\boldsymbol{x_j})}$ is maximized and

$\alpha_i y_i \boxed{\varphi(\boldsymbol{x_i})^T \varphi(\boldsymbol{x_j})} + b$

(1) $\Sigma\alpha_i y_i = 0$

(2) $0 \le \alpha_i \le C$ for all $\alpha_i$

*What if we can by-pass this explicit mapping step?*

- SVM does not need direct access to the original feature space, i.e., original data representation **x**

- It only requires access to the dot products $\mathbf{x}_i^T\mathbf{x}_j$

- The inner products

Assignment Project Exam Help

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

https://eduassistpro.github.io/

*Can be regarded as a measure of simil            data points (think cosine similarity)*

Add WeChat edu_assist_pro

Find $\alpha_1...\alpha_L$ such that
$\mathbf{Q}(\mathbf{\alpha}) = \Sigma\alpha_i - \frac{1}{2}\Sigma\sigma\alpha_i\alpha_j y_i y_j \boxed{K(\mathbf{x}_i,\mathbf{x}_j)}$ is
maximized and
(1) $\Sigma\alpha_i y_i = 0$
(2) $0 \leq \alpha_i \leq C$ for all $\alpha_i$

$$f(\mathbf{x}) = \Sigma\alpha_i y_i \boxed{K(\mathbf{x}_i,\mathbf{x})} + b$$

- What if we have a function that compute the inner product $K(\mathbf{x}_i, \mathbf{x}_j)$ directly without explicitly performing the mapping $\Phi: \mathbf{x} \rightarrow \varphi(\mathbf{x})$

Find $\alpha_1 ... \alpha_L$ such that

Assignment Project Exam Help

$Q(\boldsymbol{\alpha}) = \Sigma\alpha_i - \frac{1}{2}\Sigma\sigma\alpha_i\alpha_jy_iy_j$

maximized and

https://eduassistpro.github.io/

$= \Sigma\alpha_iy_i\,K(\mathbf{x}_i, \mathbf{x}) + b$

(1)  $\Sigma\alpha_iy_i = 0$

(2)  $0 \leq \alpha_i \leq C$ for all $\alpha_i$

Add WeChat edu_assist_pro

- A *kernel function* is a function that is equivalent to an inner product in some feature space.

- Thus, a kernel function *implicitly* maps data to a high-dimensional space (without the need to compute each $\boldsymbol{\varphi}(\mathbf{x})$ explicitly).

- Why implicit mappings?

  - Save computation

  - The target space

- 2-dimensional vectors $\mathbf{x}=[x_1 \ x_2]$
- Let: $K(\mathbf{x}_i,\mathbf{x}_j) = (1 + \mathbf{x}_i^\mathsf{T}\mathbf{x}_j)^2$
- What mapping is this?
- Need to show that $K(\mathbf{x}_i,\mathbf{x}_j) = \boldsymbol{\varphi}(\mathbf{x}_i)^\mathsf{T}\boldsymbol{\varphi}(\mathbf{x}_j)$ for some $\boldsymbol{\varphi}$

$$K(\mathbf{x}_i,\mathbf{x}_j) = (1 + \mathbf{x}_i^\mathsf{T}\mathbf{x}_j)^2 = 1 + x_{i1}^2 x_{j1}^2 + 2 x_{i1}x_{j1}x_{i2}x_{j2} + x_{i2}^2 x_{j2}^2 + 2x_{i1}x_{j1} + 2x_{i2}x_{j2}$$
$$= [1 \ \ x_{i1}^2 \ \ \sqrt{2}\ x_{i1}x_{i2} \ \ x_{i2}^2 \ \ \sqrt{2}x_{i1} \ \ \sqrt{2}x_{i2}]^\mathsf{T} [1 \ \ x_{j1}^2 \ \ \sqrt{2}\ x_{j1}x_{j2} \ \ x_{j2}^2 \ \ \sqrt{2}x_{j1} \ \ \sqrt{2}x_{j2}]$$
$$= \boldsymbol{\varphi}(\mathbf{x}_i)^\mathsf{T}\boldsymbol{\varphi}(\mathbf{x}_j),$$

where $\boldsymbol{\varphi}(\mathbf{x}) = [1 \ \ x_1^2 \ \ \sqrt{2}\ x_1 x_2 \ \ x_2^2 \ \ \sqrt{2}x_1 \ \ \sqrt{2}x_2]$

- Not all 'similarity' measures are proper kernels

$$K(\mathbf{x}_i,\mathbf{x}_j)=(1 + \mathbf{x}_i^\mathsf{T}\mathbf{x}_j)^2$$

$$K(\mathbf{x}_i,\mathbf{x}_j)=(1 + \mathbf{x}_i^\mathsf{T}\mathbf{x}_j)^3$$

- For some functions $K$ ... $\phi(\mathbf{x}_i)^\mathsf{T}\phi(\mathbf{x}_j)$ can be cumbersome.

- Mercer's theorem:

  *Every positive semi-definite symmetric function is a kernel*

- Positive semi-definite symmetric functions correspond to a positive semi-definite symmetric Gram matrix:

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

$$
K = \begin{array}{|c|c|c|c|c|}
\hline
K(\mathbf{x}_1,\mathbf{x}_1) & K(\mathbf{x}_1,\mathbf{x}_2) & K(\mathbf{x}_1,\mathbf{x}_3) & & K(\mathbf{x}_1,\mathbf{x}_n) \\
\hline
K(\mathbf{x}_2,\mathbf{x}_1) & K(\mathbf{x}_2,\mathbf{x}_2) & K(\mathbf{x}_2,\mathbf{x}_3) & & K(\mathbf{x}_2,\mathbf{x}_n) \\
\hline
 & & & & \\
\hline
\ldots & \ldots & \ldots & \ldots & \ldots \\
\hline
K(\mathbf{x}_n,\mathbf{x}_1) & K(\mathbf{x}_n,\mathbf{x}_2) & K(\mathbf{x}_n,\mathbf{x}_3) & \ldots & K(\mathbf{x}_n,\mathbf{x}_n) \\
\hline
\end{array}
$$

Non examinable

- Linear:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\mathsf{T} \mathbf{x}_j$$

  – Mapping Φ:   $\mathbf{x} \rightarrow \boldsymbol{\varphi}(\mathbf{x})$, where $\boldsymbol{\varphi}(\mathbf{x})$ is $\mathbf{x}$ itself

Assignment Project Exam Help

- Polynomial of power https://eduassistpro.github.io/

  $i$   $j$

  Add WeChat edu_assist_pro

  – Mapping Φ:   $\mathbf{x} \rightarrow \boldsymbol{\varphi}(\mathbf{x})$, where $\boldsymbol{\varphi}$ $\binom{}{p}$ dimensions

- Gaussian (Radial-Basis Function (RBF)):

$$K(\mathbf{x}_i, \mathbf{x}_j) = e^{\frac{-\left\|x_i - x_j\right\|^2}{2\sigma^2}}$$

Assignment Project Exam Help

- Mapping Φ: **x**  https://eduassistpro.github.io/ -dimensional: every point
  is mapped to *a f*

Add WeChat edu_assist_pro

- Dual problem formulation:

    Find $\alpha_1...\alpha_L$ such that

    $\mathbf{Q(\alpha)} = \Sigma\alpha_i - \frac{1}{2}\Sigma\Sigma\alpha_i\alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$ is maximized and

    (1) $\Sigma\alpha_i y_i = 0$

    (2) $C \geq \alpha_i \geq 0$

- The classifier function is:

    $$f(\mathbf{x}) = \Sigma\alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b$$

- Optimization techniques for finding $\alpha_i$'s remain the same!

- Are we guaranteed that the kernel trick will make the data linearly separable?
  - No
  - But usually work

Assignment Project Exam Help

- How to find the suitabl meters?

https://eduassistpro.github.io/

  - Method: Using  M- te

Add WeChat edu_assist_pro

- SVM is inherently a binary classifier

- Extension to multiclass:

  - One-versus-all: build M classifiers for M classes. Choose class with largest margin for test data

  - One-versus-one: one classifier per pair of classes (M(M-1)/2 classifiers in total                                        ost classifiers

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

- SVMs were originally proposed by Boser, Guyon and Vapnik in 1992 and gained increasing popularity in late 1990s.

- SVMs are currently among the best performers for a number of classification tasks ranging from text to genomic data.

- SVMs can be applied to complex data types beyond feature vectors (e.g. graphs, sequences, relational data) by designing kernel functions for such data.

- SVM techniques hav                                                                 r of tasks such as regression [Vapnik et Al. '97], principal                    analysis [Schölkopf et al. '99], etc.

- Most popular optimization algorithms for SVMs use *decomposition* to hill-climb over a subset of $\alpha_i$'s at a time, e.g. SMO [Platt '99] and  [Joachims '99]

-  Tuning SVMs remains a black art:  selecting a specific kernel and parameters is usually done in a try-and-see manner.

- [1]https://static1.squarespace.com/static/58851af9ebbd1a30e98fb283/t/58902fbae4fcb5398aeb7505/1485844411772/SVM+Explained.pdf

- [2] A Tutorial on Support Vector Machines for Pattern Recognition

- [3] Demo: http://cs.stanford.edu/People/karpathy/svmjs/demo/

  - (Note: C is the i

- [4]  Demo: http://ww VMDemo.zip

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

- What is the intuition of Support Vector Machines (SVMs)?
- How to formulate and solve SVM?
- What is linear and non-linear SVM?

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro