



Outlines

Overview
Introduction
Linear Algebra
Probability
Linear Regression
Linear Regression
Linear Classification 1
Linear Classification 2
Kernel Methods
Sparse Kernel Methods
Structure Models and EM 1
Structure Models and EM 2
Neural Networks 1
Neural Networks 2
Deep Component Analysis
Generators
Graphical Models 1
Graphical Models 2
Graphical Models 3
Sampling
Sequential Data 1
Sequential Data 2

Statistical Machine Learning

Assignment Project Exam Help

Christian Walder

<https://eduassistpro.github.io>

College of Engineering and Computer Science
The Australian National University

Add WeChat edu_assist_pro

Canberra

Semester One, 2020.

(Many figures from C. M. Bishop, "Pattern Recognition and Machine Learning")



Classification

Generalised Linear
Model

Discriminant Functions

Feature Engineering

Discriminant Analysis

The Perceptron
Algorithm

Assignment Project Exam Help

Part V

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr



Classification

Generalised Linear
Model

Discriminant Functions

Perceptron
Algorithm

The Perceptron
Algorithm

Assignment Project Exam Help

- Estimate best predictor = training = learning

Giv

1

2

3

4

5

6

7

<https://eduassistpro.github.io>

Calculate the optimal parameter (w)

Model uncertainty using the Bayesian approach

Implement and compute (the algorithm in python)

Interpret and diagnose results

Add WeChat edu_assist_pro



- Goal : Given input data \mathbf{x} , assign it to one of K discrete classes \mathcal{C}_k where $k = 1, \dots, K$.
- Divide the input space into different regions.
- Eq

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Length of petal [in cm] vs sepal [cm] for three types of flowers
(Iris Setosa, Iris Versicolor, Iris Virginica).

How to represent binary class labels?



Classification

Generalised Linear
Model

Discriminant Functions

Perceptron
Discriminant

The Perceptron
Algorithm

Assignment Project Exam Help

- Class labels are no longer real values as in regression, but a discrete

- Two
(t)

- Can

with only two values possible for the probability

- Note: Other conventions to map classes into integers possible, check the setup.

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

How to represent multi-class labels?



Assignment Project Exam Help

- If there are more than two classes ($K > 2$), we call it a multi-class setup.

- Oft
len

- Exa

class C_2 will be encoded as the target vector

<https://eduassistpro.github.io>

Add WeChat [edu_assist_pro](https://eduassistpro.github.io)

- Note: Other conventions to map multi-classes are possible, check the setup.

$$\mathbf{t} = (0, 1, 0, 0)^T$$



- Idea: Use again a Linear Model as in regression: $y(\mathbf{x}, \mathbf{w})$ is a linear function of the parameters \mathbf{w}

$$y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^\top \phi(\mathbf{x})$$

- But
Exam

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr



- Apply a mapping $f : \mathbb{R} \rightarrow \mathbb{Z}$ to the linear model to get the discrete class labels.

Assignment Project Exam Help

$$y(\mathbf{x}_n, \mathbf{w}) = f(\mathbf{w}^\top \phi(\mathbf{x}_n))$$

- Acti
- Lin

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

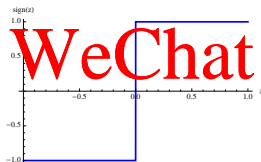


Figure: Example of an activation function $f(z) = \text{sign}(z)$.

Three Models for Decision Problems



Classification

Generalised Linear
Model

Discriminant Functions

Bayes' theorem
Discriminant

The Perceptron
Algorithm

- Find a **discriminant function** $f(\mathbf{x})$ which maps each input directly onto a class label
- Discriminative Models

1

2

<https://eduassistpro.github.io>

- Generative Models

1

Solve the inference problem of determining the class conditional probabilities $p(\mathbf{x} | C_k)$

2

Also infer the prior class probabilities

3

Use Bayes' theorem to find the posterior

4

Alternatively model the joint distribution $p(\mathbf{x}, C_k)$ directly.

5

Use decision theory to assign each new \mathbf{x} to one of the classes.

Add WeChat edu_assist_pro



Definition

A **discriminant** is a function that maps from an input vector \mathbf{x} to one of K classes, denoted by k .

- Co
- Co

$$y(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + w_0$$

such that \mathbf{x} being assigned to class \mathcal{C}_1
class \mathcal{C}_2 otherwise.

- **weight vector** \mathbf{w}
- **bias** w_0 (sometimes $-w_0$ called **threshold**)

Classification

Generalised Linear
Model

Discriminant Functions

Perceptron
Discriminant

The Perceptron
Algorithm



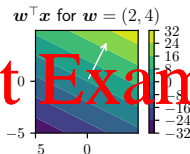
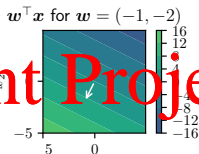
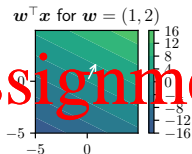
Classification

Generalised Linear Model

Discriminant Functions

Perceptron Algorithm

The Perceptron Algorithm



- Graph of the linear function $w^T x$ is a hyper-plane.
- The hyper-plane $w^T x = 0$ is the decision boundary.
- Projecting x on that hyper-plane means find $\arg \min_x \|x - x_\perp\|$ subject to the constraint $w^T x_\perp + w_0 = 0$. Geometrically: move in the direction w until you reach the hyper-plane.
- Rate of change of function value in that direction is

$$\frac{d}{da} \left(a \frac{w^T x}{\|w\|} \right) \frac{w}{\|w\|} = a \|w\|.$$

- The length $\left\| a \frac{w}{\|w\|} \right\| = \frac{a}{\|w\|} \|w\| = a$.
- For a fixed change in $w^T x$, $a \propto \frac{1}{\|w\|}$.



Classification

Generalised Linear
Model

Discriminant Functions

Perceptron
Algorithm

The Perceptron
Algorithm

Assignment Project Exam Help

- Decision boundary $y(\mathbf{x}) = 0$ is a $(D - 1)$ -dimensional
hyp
surf

- \mathbf{w} is
- Pro $\mathbf{w}^T \mathbf{x}_A \geq \mathbf{w}^T \mathbf{x}_B$ decision surface. Then,

Add WeChat edu_assist_pr

Two Classes

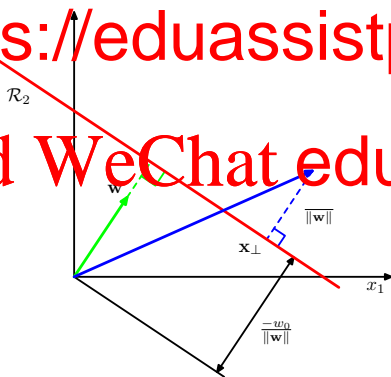
- $y(\mathbf{x})$ gives a **signed** measure of the perpendicular distance r **from** the decision surface **to** \mathbf{x} , that is $r = y(\mathbf{x})/\|\mathbf{w}\|$.

Assignment Project Exam Help

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|^2} + w_0 = r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|^2} + \mathbf{w}^T \mathbf{x} + w_0 = r \frac{\|\mathbf{w}\|^2}{\|\mathbf{w}\|^2} + \mathbf{w}^T \mathbf{x} + w_0 = r + \mathbf{w}^T \mathbf{x} + w_0$$

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro





Classification

Generalised Linear
Model

Discriminant Functions

Perceptron Algorithm

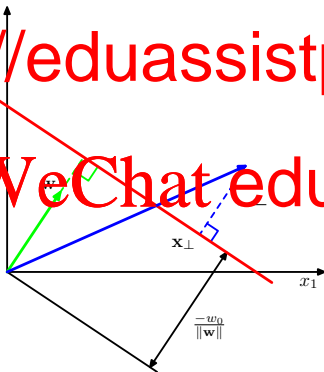
The Perceptron
Algorithm

- The normal distance from the origin to the decision surface is therefore

$$-\frac{y(\mathbf{0})}{\|\mathbf{w}\|} = -\frac{w_0}{\|\mathbf{w}\|}$$

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro





Classification

Generalised Linear
Model

Discriminant Functions

Perceptron
Discriminant

The Perceptron
Algorithm

Assignment Project Exam Help

- More compact notation : Add an extra dimension to the inp
- Als

<https://eduassistpro.github.io>

(if it helps, you may think of $\tilde{\mathbf{w}}^T$ as a func

- Decision surface is now a D -dimension
 $D + 1$ -dimensional expanded input space.

Add WeChat edu_assist_pr

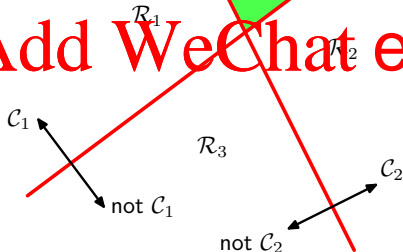


- Number of classes $K > 2$
- Can we combine a number of two-class discriminant functions using $K - 1$ one-versus-the-rest classifiers?

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



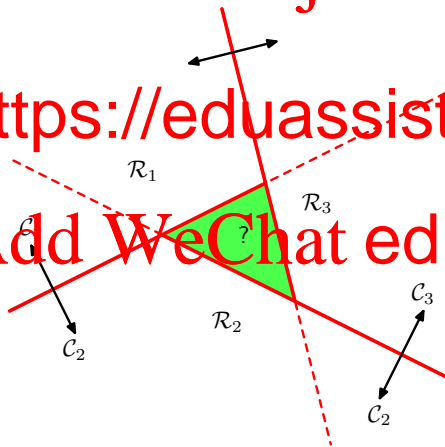


- Number of classes $K > 2$
- Can we combine a number of two-class discriminant functions using $K(K-1)/2$ one-versus-one classifiers?

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro





Classification

Generalised Linear
Model

Discriminant Functions

Perceptron

Discriminant

The Perceptron
Algorithm

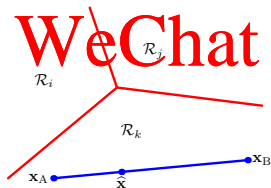
- Number of classes $K > 2$
- Solution: Use K linear functions

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

- Ass
 - De
- $y_k(\mathbf{x})$ j

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro





Classification

Generalised Linear
Model

Discriminant Functions

Perceptron Algorithm

Discriminant

The Perceptron
Algorithm

Assignment Project Exam Help

- Regression with a linear function of the model parameters and minimisation of sum-of-squares error function resulted in a classifier
- Is this a linear classifier?
- Given a set of training data, how can we find the model parameters?
- Use 1-of- K binary coding scheme.
- Each class is described by its own linear model

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

$$y_k(\mathbf{x}) = \mathbf{w}_k^\top \mathbf{x} + w_{k0} \quad k = 1, \dots, K$$



Classification

Generalised Linear
Model

Discriminant Functions

Perceptron
Discriminant

The Perceptron
Algorithm

- With the conventions

Assignment Project Exam Help

$$\mathbf{w} = \begin{bmatrix} w_0 \\ \mathbf{w}_k \end{bmatrix} \in \mathbb{R}^{D+1}$$

<https://eduassistpro.github.io>

- we get for the (vector valued) discriminant func

Add WeChat edu_assist_pro

(if it helps, you may think of $\tilde{\mathbf{W}}^\top$ as a **vector-valued** function).

- For a new input \mathbf{x} , the class is then defined by the index of the largest value in the row vector $\mathbf{y}(\mathbf{x})$

Determine $\tilde{\mathbf{W}}$



- Given a training set $\{\mathbf{x}_n, \mathbf{t}_n\}$ where $n = 1, \dots, N$, and \mathbf{t}_n is the class in the 1-of-K coding scheme.
- Define a matrix \mathbf{T} where row n corresponds to \mathbf{t}_n^\top .
- The

(ch

<https://eduassistpro.github.io>

$$E_D(\tilde{\mathbf{W}}) = \frac{1}{2} \text{tr} \left(\tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top - \tilde{\mathbf{X}} \mathbf{T} \tilde{\mathbf{X}}^\top \right)$$

- The minimum of $E_D(\tilde{\mathbf{W}})$ will be reached for

$$\tilde{\mathbf{W}} = (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^\top \mathbf{T} = \tilde{\mathbf{X}}^\dagger \mathbf{T}$$

where $\tilde{\mathbf{X}}^\dagger$ is the pseudo-inverse of $\tilde{\mathbf{X}}$.

Add WeChat edu_assist_pro



- The discriminant function $\mathbf{y}(\mathbf{x})$ is therefore

$$\mathbf{y}(\mathbf{x}) = \mathbf{W}^T \mathbf{x} = \mathbf{J}^T (\mathbf{X}^T)^T \mathbf{x}$$

wh
inp

- Inten
con
also obey the same constraint

<https://eduassistpro.github.io>
Add WeChat edu_assist_pro

- For the 1-of- K coding scheme, the sum of all components in \mathbf{t}_n is one, and therefore all components of $\mathbf{y}(\mathbf{x})$ will sum to one. BUT: the components are not probabilities, as they are not constraint to the interval $(0, 1)$.

Deficiencies of the Least Squares Approach

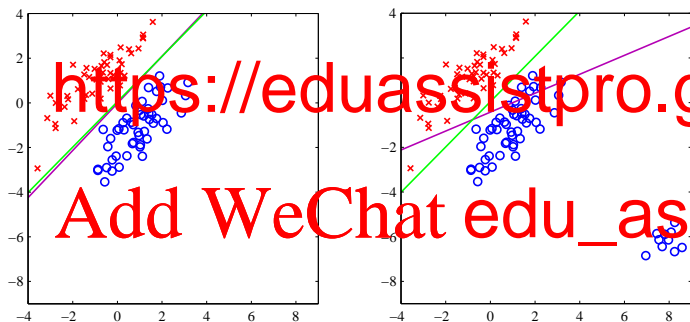


Magenta curve :

Decision Boundary for the least squares approach

Green curve :

Decision boundary for the logistic regression (described later)



(Imagine heat-maps of the quadratic penalty function, similarly to those of the linear functions earlier in the slides.)

Classification

Generalised Linear
Model

Discriminant Functions

Perceptron
Discriminant

The Perceptron
Algorithm

Deficiencies of the Least Squares Approach



Classification

Generalised Linear
Model

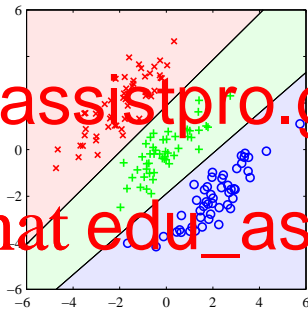
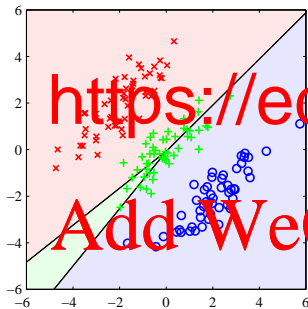
Discriminant Functions

Perceptron
Discriminant

The Perceptron
Algorithm

Left plot : Decision Boundary for least squares

Right plot : Boundaries for logistic regression (described later)



<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



Classification

Generalised Linear
Model

Discriminant Functions

Fisher's Linear
Discriminant

The Perceptron
Algorithm

Assignment Project Exam Help

- View linear classification as dimensionality reduction.

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

If y

- But
spa

- Projection always means loss of information.
- For classification we want to preserve the class in one dimension.
- Can we find a projection which maximally preserves the class separation ?

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Fisher's Linear Discriminant

Statistical Machine
Learning

© 2020

Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University



Classification

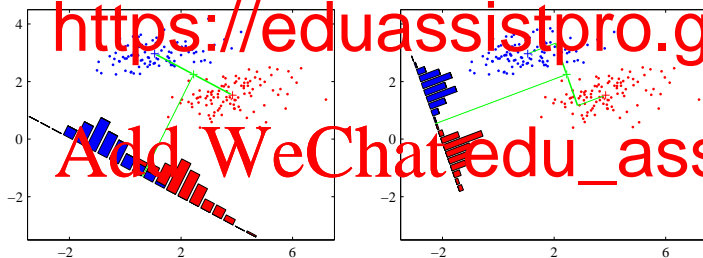
Generalised Linear
Model

Discriminant Functions

Fisher's Linear
Discriminant

The Perceptron
Algorithm

Samples from two classes in a two-dimensional input space and their histogram when projected to two different one-dimensional spaces.



<https://eduassistpro.github.io>
Add WeChat: edu_assist_pr

Fisher's Linear Discriminant - First Try

Statistical Machine
Learning

© 2020

Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University



Classification

Generalised Linear
Model

Discriminant Functions

Fisher's Linear
Discriminant

The Perceptron
Algorithm

- Given N_1 input data of class C_1 , and N_2 input data of class C_2 , calculate the centres of the two classes

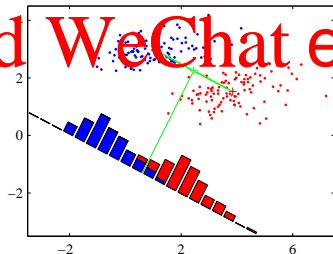
$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in C_1} \mathbf{x}_n, \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in C_2} \mathbf{x}_n$$

- Ch
proj

<https://eduassistpro.github.io>

- Problem with non-uniform covariance

Add WeChat edu_assist_pro





- Measure also the within-class variance for each class

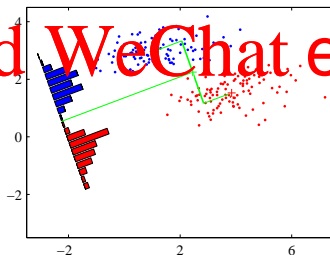
$$s_k^2 = \sum_{y_n \in \mathcal{C}_k} (y_n - m_k)^2$$

where $y_n = \mathbf{w}^\top \mathbf{x}_n$.

- Ma

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr



Primer: Bilinear form with a Covariance Matrix

Let

$$\boldsymbol{\mu} = \mathbb{E}[\mathbf{x}]$$

$$\boldsymbol{\Sigma} = \text{cov}[\mathbf{x}]$$

$$= \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top].$$

Then

<https://eduassistpro.github.io>

$$= \mathbb{E}[(\mathbf{w}^\top \mathbf{x} - \mathbf{w}^\top \mathbb{E}[\mathbf{x}])^2]$$

$$= \mathbb{E}[(\mathbf{w}^\top \mathbf{x} - \mathbf{w}^\top \boldsymbol{\mu})^2]$$

$$= \mathbb{E}[(\mathbf{w}^\top \mathbf{x} - \mathbf{w}^\top \boldsymbol{\mu})(\mathbf{w}^\top \mathbf{x} - \mathbf{w}^\top \boldsymbol{\mu})]$$

$$= \mathbb{E}[(\mathbf{w}^\top \mathbf{x} - \mathbf{w}^\top \boldsymbol{\mu})(\mathbf{x}^\top \mathbf{w} - \boldsymbol{\mu}^\top \mathbf{w})]$$

$$= \mathbb{E}[\mathbf{w}^\top (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{w}]$$

$$= \mathbf{w}^\top \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top] \mathbf{w}$$

$$= \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w}.$$





- The Fisher criterion can be rewritten as

$$J(\mathbf{w}) = \frac{\mathbf{w}^\top \mathbf{S}_B \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_W \mathbf{w}}$$

- \mathbf{S}_B is the **between-class** covariance

so b

the variance of the projection of the means

- \mathbf{S}_W is the **within-class** covariance

$$\mathbf{S}_W = \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T$$

so so by the previous slide and

$\mathbf{w}^\top (A + B) \mathbf{w} = \mathbf{w}^\top A \mathbf{w} + \mathbf{w}^\top B \mathbf{w}$, the denominator of $J(\mathbf{w})$ is:
(the variance of the projection of the points in class \mathcal{C}_1) +
(the variance of the projection of the points in class \mathcal{C}_2)



Assignment Project Exam Help

- The Fisher criterion

$$\underline{\mathbf{w}^\top \mathbf{S}_B \mathbf{w}}$$

has

<https://eduassistpro.github.io>

$$\mathbf{w} \propto \mathbf{S}_W^{-1}(\mathbf{m}_2 - \mathbf{m}_1)$$

- Fisher's Linear discriminant is NOT a discriminant. It can be used to construct one by choosing a threshold projection space.

Add WeChat edu_assist_pro

Fisher's Discriminant For Multi-Class



- Assume that the dimensionality of the input space D is greater than the number of classes K .
- Use $D' > 1$ linear 'features' $y_k = \mathbf{w}^\top \mathbf{x}$ and write everything in vector form (with no bias term)

$$\mathbf{y} = \mathbf{W}^\top \mathbf{x}.$$

- The cov

$$\mathbf{S}_W = \sum_{k=1}^K \mathbf{S}_k$$

where

$$\mathbf{S}_k = \sum_{n \in \mathcal{C}_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^\top$$
$$\mathbf{m}_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k} \mathbf{x}_n$$



- Between-class covariance

$$\mathbf{S}_B = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{n})(\mathbf{m}_k - \mathbf{n})^T$$

wh

<https://eduassistpro.github.io>

$n=1$

- One possible way to define a function of when the between-class covariance is large a within-class covariance is small is given by

$$J(\mathbf{W}) = \text{tr} \{ (\mathbf{W}^\top \mathbf{S}_W \mathbf{W})^{-1} (\mathbf{W}^\top \mathbf{S}_B \mathbf{W}) \}$$

- The maximum of $J(\mathbf{W})$ is determined by the D' eigenvectors of $\mathbf{S}_W^{-1} \mathbf{S}_B$ with the largest eigenvalues.

Add WeChat edu_assist_pro



Classification

Generalised Linear
Model

Discriminant Functions

Perceptron
Algorithm

The Perceptron
Algorithm

- Frank Rosenblatt (1928 - 1969)
- "Principles of neurodynamics: Perceptions and the theory of brain mechanisms" (Spartan Books, 1962)

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr



Classification

Generalised Linear
Model

Discriminant Functions

Perceptron

Discriminant

The Perceptron
Algorithm

• Perceptron ("MARK 1") was the first computer which could learn new skills by trial and error

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

The Perceptron Algorithm



Classification

Generalised Linear
Model

Discriminant Functions

Perceptron
Algorithm

The Perceptron
Algorithm

- Two class model
- Create feature vector $\phi(\mathbf{x})$ by a fixed nonlinear transformation of the input \mathbf{x} .
- Generalised linear model

wit

- non

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

- Target coding for perceptron

$$t = \begin{cases} +1, & \text{if } C_1 \\ -1, & \text{if } C_2 \end{cases}$$



Classification

Generalised Linear
Model

Discriminant Functions

Perceptron
Discriminant

The Perceptron
Algorithm

Assignment Project Exam Help

- Idea: Minimise total number of misclassified patterns.
- Problem: As a function of \mathbf{w} , this is piecewise constant and t
- Bet wa
- **Perceptron Criterion**: Add the errors for all patterns belonging to the set of misclassified patterns

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

$$E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^T \phi(\mathbf{x}_n)$$



- Perceptron Criterion (with notation $\phi_n = \phi(\mathbf{x}_n)$)

$$E_P(\mathbf{w}) = \sum_{n \in \mathcal{M}} y_n^T \phi_n t_n \\ \equiv E_P^{(n)}(\mathbf{w})$$

- On

1 <https://eduassistpro.github.io>

- 2 Update the weight vector \mathbf{w} by

Add $\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E$
where

$$\nabla E_P^{(n)}(\mathbf{w}) = \begin{cases} -\phi_n t_n & \text{if } (\mathbf{w}^{(\tau)})^\top \phi(\mathbf{x}_n) \cdot t_n \leq 0 \\ 0 & \text{otherwise.} \end{cases}$$

- As $y(\mathbf{x}, \mathbf{w})$ is invariant to the norm of \mathbf{w} , we may set $\eta = 1$.

The Perceptron Algorithm - Update 1

Statistical Machine
Learning

© 2020
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University



Classification

Generalised Linear
Model

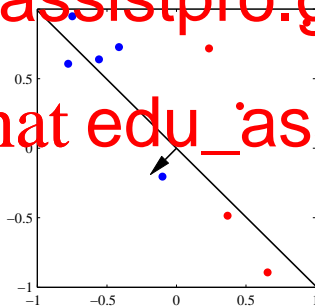
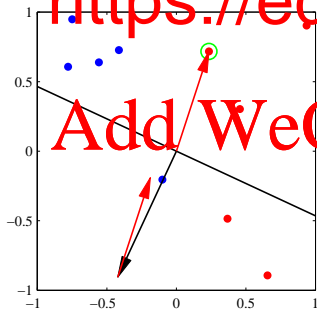
Discriminant Functions

Perceptron Learning
Algorithm

The Perceptron
Algorithm

Update of the perceptron weights from a misclassified pattern

(green)
 $\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \varphi_n t_n$



The Perceptron Algorithm - Update 2

Statistical Machine
Learning

© 2020
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University



Classification

Generalised Linear
Model

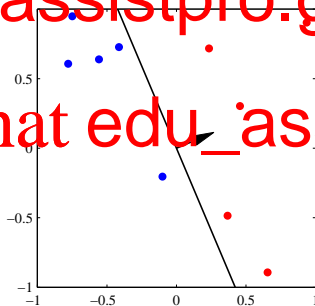
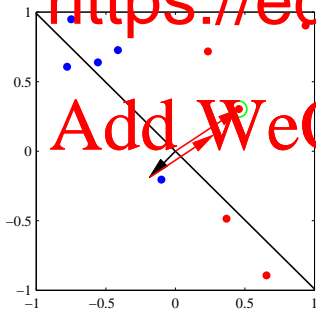
Discriminant Functions

Perceptron Algorithm

The Perceptron
Algorithm

Update of the perceptron weights from a misclassified pattern

(green)
 $\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \varphi_n t_n$



The Perceptron Algorithm - Convergence



- Does the algorithm converge ?

- For a single update step, letting $\eta = 1$, and considering the error from a single point,

— $(\tau+1)T$

$(\tau)T$

$(\tau)T$

bec
gra
function.

- BUT: contributions to the error from the other mi
patterns might have increased.

- AND: some correctly classified patterns might
misclassified.

- Perceptron Convergence Theorem** : If the training set is
linearly separable, the perceptron algorithm is guaranteed
to find a solution in a finite number of steps.