

# Assignment Project Exam Help

Supervised Learning – Classification

<https://eduassistpro.github.io>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

Last Revision 11 March 2023

## Acknowledgements

Material derived from slides for the book

"Elements of Statistical Learning (2nd Ed.)" by T. Hastie,  
R. Tibshirani & J. Friedman. Springer (2009)

<http://statweb.stanford.edu/~tibs/ElemStatLearn>

Material derived from slides for the book

"Machine Learning: A Probabilistic Perspective" by P. Murphy

MIT Press (2012)

<http://www.cs.cmu.edu/~pemurphy/mle.html>

Material derived from slides for the book

"Machine Learning"

Cambridge University Press

<http://www.cs.cmu.edu/~pemurphy/ml.html>

Material derived from slides for the book

"Bayesian Reasoning and Machine Learning" by D. Barber

Cambridge University Press (2012)

<http://www.cs.cmu.edu/~d-barber/ml.html>

Material derived from slides for the book

"Machine Learning" by T. Mitchell

McGraw-Hill (1997)

<http://www-2.cs.cmu.edu/~tom/mlbook.html>

Material derived from slides for the course

"Machine Learning" by A. Srinivasan

BITS Pilani, Goa, India (2016)

# Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

## Aims

This lecture will introduce you to machine learning approaches to the problem of *classification*. Following it you should be able to reproduce theoretical results, outline algorithmic techniques and describe practical applications for the topics:

- des
- outl
- defi
- define the Bayes optimal classification rule in ter
- outline the Naive Bayes classification algorithm
- describe typical applications of Naive Bayes for t
- outline the Perceptron classification algorithm
- outline the logistic regression classification algorithm

Note: slides with titles marked \* are for background only.

# Introduction

Classification (sometimes called *concept learning*) methods dominate machine learning:

# Assignment Project Exam Help

... however, they often don't have convenient mathematical properties like regression

*classifier*

one of a set of

<https://eduassistpro.github.io>

We will mostly focus on their advantages and disadvantages first, and point to unifying ideas and approaches applicable. In this lecture we focus on classification methods, essentially *linear models* ...

and in later lectures we will see other, more expressive, classifier learning methods.

# Nearest neighbour classification

- Related to the simplest form of learning: rote learning or memorization

Training instances are searched for instance that most closely **resembles** new or *query* instance

⋮

⋮

- The method is *instance-based learning*; beyond simple memorization
- Intuitive idea — instances “close by”, i.e., neighbours, should be classified similarly
- Instance-based learning is *lazy learning*
- Methods: *nearest-neighbour*, *k-nearest-neighbour*, *lowess*, ...
- Ideas also important for *unsupervised* methods, e.g., clustering (later lectures)

Add WeChat `edu_assist_pro`

## Minkowski distance

# Assignment Project Exam Help

*Minkowski distance* if  $\mathcal{X} = \mathbb{R}^d$ , the *Minkowski distance* of order  $p > 0$  is defined as

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

$j=1$

where  $\|\mathbf{z}\|_p = \left( \sum_{j=1}^d |z_j|^p \right)^{1/p}$  is the  $p$ -norm (L<sub>p</sub> norm) of the vector  $\mathbf{z}$ .

Minkowski distance

• The 2-norm refers to the familiar *Euclidean distance*  
 $d(\mathbf{x}, \mathbf{y})$

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

whi

- The 1-norm denotes *Manhattan distance*:

Add WeChat [edu\\_assist\\_pro.com](https://edu_assist_pro.com)

$$\text{Dis}_1(\mathbf{x}, \mathbf{y}) = \sum_{j=1}^d |x_j - y_j|$$

This is the distance if we can only travel along coordinate axes.

## Minkowski distance

- If we now let  $p$  grow larger, the distance will be more and more dominated by the largest coordinate wise distance, from which we can infer that  $\text{Dis}_{\infty}(\mathbf{x}, \mathbf{y}) = \max_j |x_j - y_j|$ ; this is also called Chebyshev dist
- You can corr <https://eduassistpro.github.io>

vectors  $\mathbf{x}$  and  $\mathbf{y}$  differ. This is not strictly a Mi

however, we can define it as

Add WeChat edu\_assist\_pr

$$\text{Dis}_0(\mathbf{x}, \mathbf{y}) = \sum_{j=1}^d (x_j - y_j)^0 = \sum_{j=1}^d I[x_j = y_j]$$

under the understanding that  $x^0 = 0$  for  $x = 0$  and 1 otherwise.

Minkowski distance

# Assignment Project Exam Help

Sometimes the data  $\mathcal{X}$  is not naturally in  $\mathbb{R}^d$ , but if we can turn it into

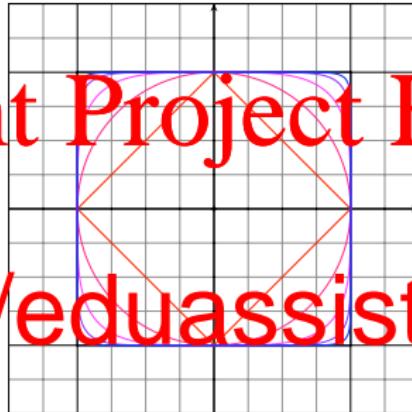
Boolean f

measure

- If  $x$  is a binary string, we can use the Hamming distance as the number of bits that need to be flipped to change  $x$  into  $y$ . Alternatively, we can see the Hamming distance as the number of bits that need to be flipped to change  $x$  into  $y$ .
- For non-binary strings of unequal length, this is called the Levenshtein edit distance or Levenshtein distance.

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

## Circles and ellipses



# Assignment Project Exam Help

<https://eduassistpro.github.io>

Lines connecting points sat order- $p$  Minkowski distance for (from inside)  $p = 0.8$ ,  $p = 1$  (Manhattan in red);  $p = 1.5$ ;  $p = 2$  (Euclidean distance, the violet circle);  $p = 4$ ;  $p = 8$ ; and  $p = \infty$  (Chebyshev distance, the blue rectangle). Notice that for points on the coordinate axes all distances agree. For the other points, our reach increases with  $p$ ; however, if we require a rotation-invariant distance metric then Euclidean distance is our only choice.

## Distance metric

# Assignment Project Exam Help

Distance metric Given an instance space  $\mathcal{X}$ , a *distance metric* is a function  $\text{Dis} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  such that for any  $x, y, z \in \mathcal{X}$ :

- $\text{dist}(x, y) \geq 0$  ;
- all of  $\text{dist}(x, y) = 0 \iff x = y$ ;
- $\text{dist}(x, y) = \text{dist}(y, x)$ ;
- detours can not shorten the distance:  
$$\text{Dis}(x, z) \leq \text{Dis}(x, y) + \text{Dis}(y, z)$$

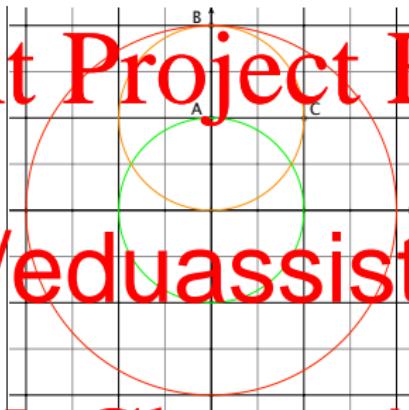
If the second condition is weakened to a non-strict inequality  $\text{Dis}(x, y) \leq \text{Dis}(x, z) + \text{Dis}(z, y)$ ,

$\text{Dis}(x, y)$  may be zero even if  $x \neq y$  – the function  $\text{Dis}$  is called a *pseudo-metric*.

The triangle inequality – Minkowski distance for  $p = 2$

Assignment Project Exam Help

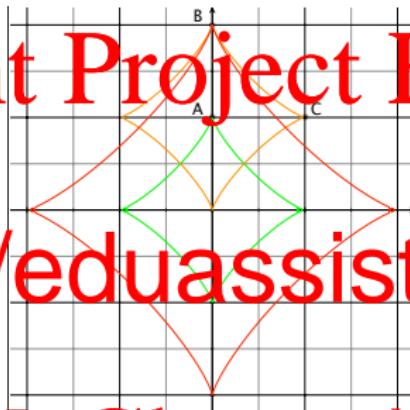
<https://eduassistpro.github.io>



Add WeChat `edu_assist_pro`

The green circle connects points the same Euclidean distance from A. The orange circle shows that B and C are equidistant from A. The red circle demonstrates that C is closer to the origin than B, which conforms to the triangle inequality.

The triangle inequality – Minkowski distance for  $p \leq 1$



<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pro

With Manhattan distance ( $p = 1$ ), B and C are equidistant from A.

and also equidistant from A. With  $p < 1$  (here,  $p = 0.8$ ) C is further away from the origin than B; since both are again equidistant from A, it follows that travelling from the origin to C via A is quicker than going there directly, which violates the triangle inequality.

## Mahalanobis distance \*

# Assignment Project Exam Help

Often, the shape of the ellipse is estimated from data as the inverse of the covariance matrix:  $\mathbf{M} = \Sigma^{-1}$ . This leads to the definition of the *Mahalanobis*

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

Using the covariance matrix in this way has the effect of de  
normalising the features

Add WeChat edu\_assist\_pro

Clearly, Euclidean distance is a special case of Mahalanobis distance with  
the identity matrix  $\mathbf{I}$  as covariance matrix:  $\text{Dis}_2(\mathbf{x}, \mathbf{y}) = \text{Dis}_M(\mathbf{x}, \mathbf{y} | \mathbf{I})$ .

## Means and distances

The arithmetic mean minimises squared Euclidean distance  
The arithmetic mean  $\mu$  of a set of data points  $D$  in a Euclidean space is the unique point that minimises the sum of squared Euclidean distances to those data

Proof. We write

denotes the

vector of partial derivatives with respect to  
the zero vector:

$$\nabla_y \sum_{x \in D} \|x - y\|^2 = -2 \sum_{x \in D} (x - y) \quad \text{at } y = 0$$

from which we derive  $y = \frac{1}{|D|} \sum_{x \in D} x = \mu$ .

## Means and distances

# Assignment Project Exam Help

Notice that minimising the sum of squared Euclidean distances of a given set of points is the same as minimising the average squared Euc

- You can it be minimised by choosing a point as exemplar?
- This point is known as the *geometric median*. It corresponds to the median or middle value of a set of numbers. However, for multivariate data there is no close approximation to the geometric median, which needs to be calculated by successive approximation.

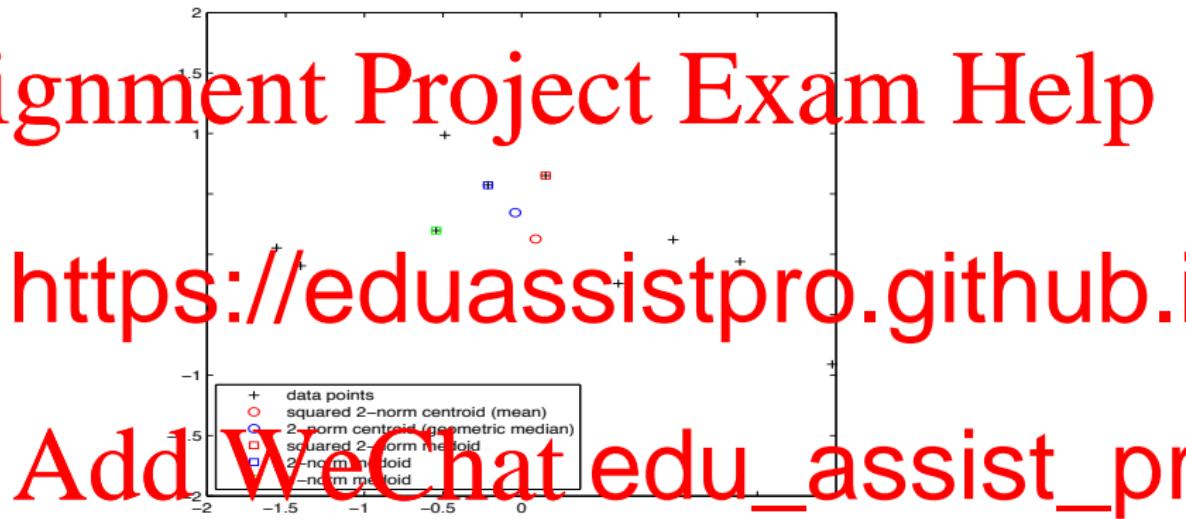
Add WeChat `edu_assist_pro`

## Means and distances

# Assignment Project Exam Help

- In certain situations it makes sense to restrict an exemplar to be one of the given data points. In that case, we speak of a *medoid*, to dist to oc e
  - Fin tota <https://eduassistpro.github.io/> that minimises it. Regardless of the distance metric, an  $O(n^2)$  operation for  $n$  points.
  - So for medoids there is no computational reason to prefer one distance metric over another.
  - There may be more than one medoid.
- Add WeChat edu\_assist\_pro

## Centroids and medoids



A small data set of 10 points, with circles indicating centroids and squares indicating medoids (the latter must be data points), for different distance metrics. Notice how the outlier on the bottom-right 'pulls' the mean away from the geometric median; as a result the corresponding medoid changes as well.

# The basic linear classifier is distance-based

- The basic linear classifier constructs the decision boundary as the perpendicular bisector of the line segment connecting the two exemplars (one for each class).

- An al

refe

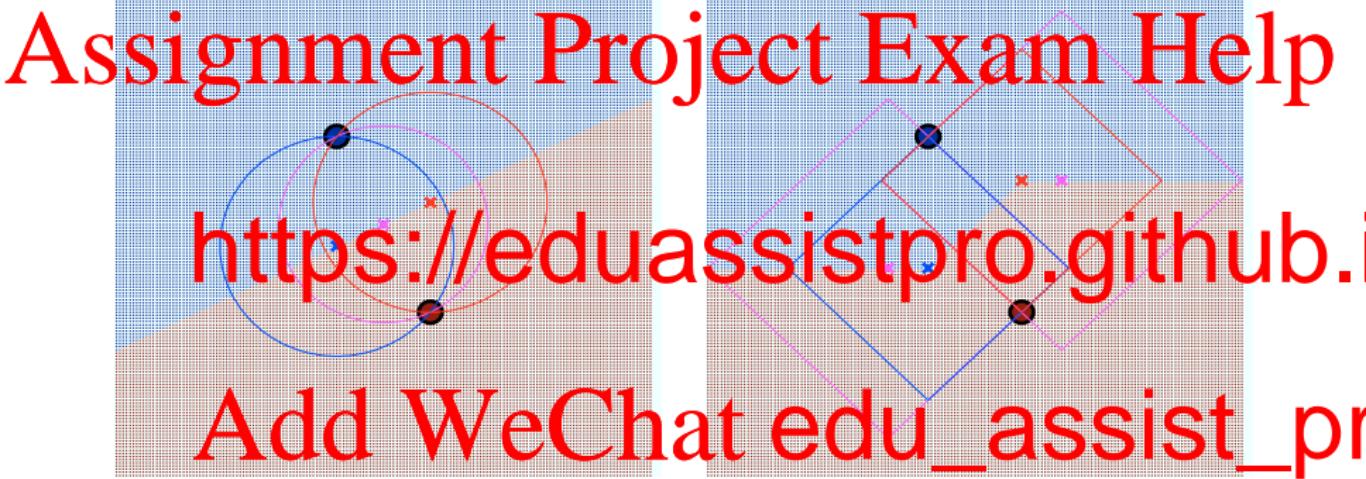
is ne

<https://eduassistpro.github.io>

equivalently, classify an instance to the class of the *nearest* exemplar.

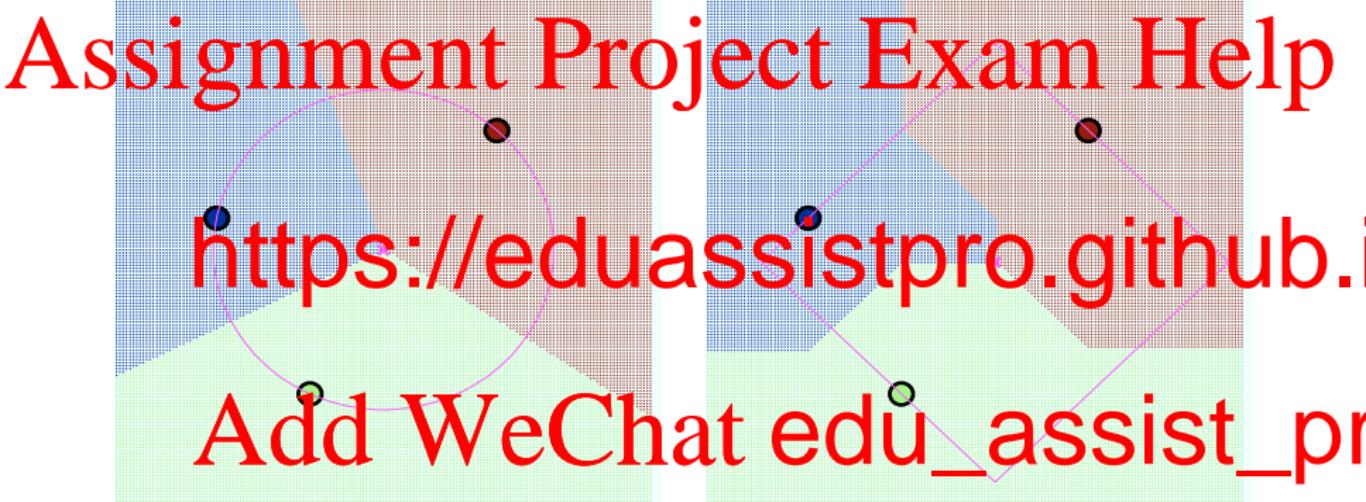
- If we use Euclidean distance as our closeness measure, geometry tells us we get exactly the same decision.
- So the basic linear classifier can be interpreted from a perspective as constructing exemplars that minimise squared Euclidean distance within each class, and then applying a nearest-exemplar decision rule.

## Two-exemplar decision boundaries



(left) For two exemplars the nearest-exemplar decision rule with Euclidean distance results in a linear decision boundary coinciding with the perpendicular bisector of the line connecting the two exemplars. (right) Using Manhattan distance the circles are replaced by diamonds.

## Three-exemplar decision boundaries



(left) Decision regions defined by the 2-norm nearest-exemplar decision rule for three exemplars. (right) With Manhattan distance the decision regions become non-convex.

## Distance-based models

# Assignment Project Exam Help

To summarise, the main ingredients of distance-based models are

- dist
  - Ma
  - exe
  - <https://eduassistpro.github.io/>
- distance metric, or medoids that find the most centrally located data point; and
- distance-based decision rules, which take a vote between nearest exemplars.
- Add WeChat edu\_assist\_pro

## Nearest Neighbour

Stores all training examples  $\langle x_i, f(x_i) \rangle$ .

Nearest neighbour:

- Given  $x_q$ , then  
estimate  $f(x_q)$

*k*-Nearest

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

- Given  $x_q$ , take vote among its  $k$  nearest neighbours (if discrete-valued target function)
- take mean of  $f$  values of  $k$  nearest neighbours

Add WeChat edu\_assist\_pro

$$\hat{f}(x_q) \leftarrow \frac{\sum_{i=1}^k f(x_i)}{k}$$

# *k*-Nearest Neighbour Algorithm

## Assignment Project Exam Help

- For each training example  $\langle x_i, f(x_i) \rangle$ , add the example to the list  $training\_examples$

Classification

- Given

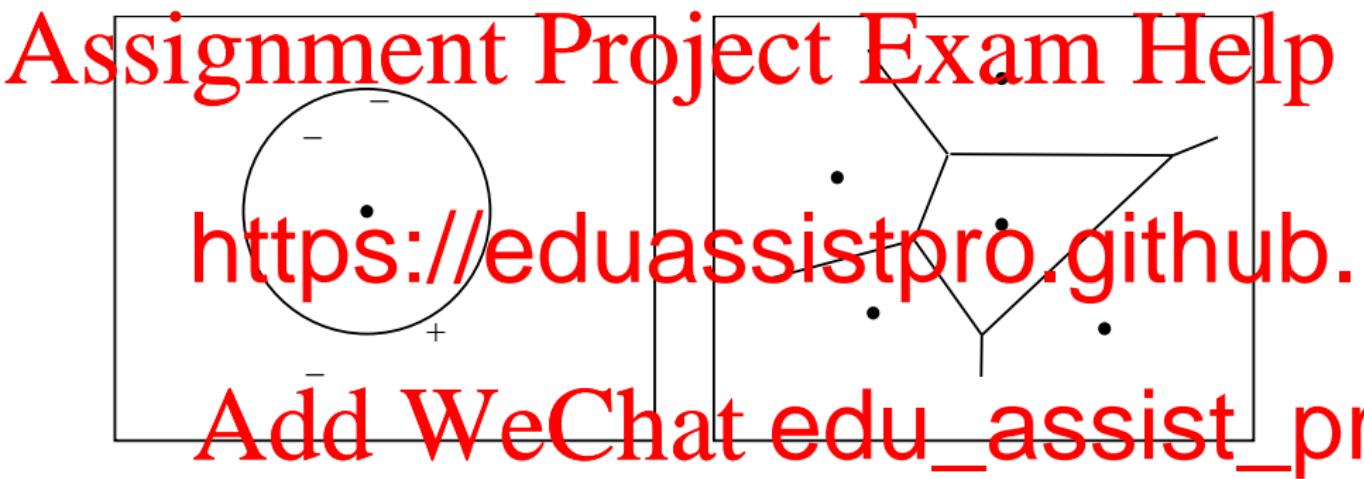
- Let  $x_1 \dots x_k$  be the  $k$  instances from *training-examples* that are nearest to  $x_q$  by the distance function

- Return

$f(x_q) \leftarrow \arg \max_{v \in V} \sum_{i=1}^k \delta(x_q, x_i)$

where  $\delta(a, b) = 1$  if  $a = b$  and 0 otherwise.

## "Hypothesis Space" for Nearest Neighbour



2 classes, + and – and query point  $x_q$ . On left, note effect of varying  $k$ .  
On right, 1–NN induces a Voronoi tessellation of the instance space.  
Formed by the perpendicular bisectors of lines between points.

## Distance function again

The distance function defines what is learned.  
Instance  $x$  is described by a feature vector (list of attribute-value pairs)

where  $a_r$   
Most com

- distance between two instances  $x_i$  an

Add WeChat edu\_assist\_pr

$$d(x_i, x_j) = \sqrt{\sum_{r=1}^m (a_r(x_i) - a_r(x_j))^2}$$

Distance function again

# Assignment Project Exam Help

Many other distance functions could be used ...

- e.g.  
diff

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

$$d(x_i, x_j) = \sqrt{\sum_{r=1}^n |a_r(x_i) - a_r(x_j)|^2}$$

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

Vector-based formalization use norm  $L_1$ ,

The idea of distance functions will appear again in *kernel methods*.

## Normalization and other issues

# Assignment Project Exam Help

- Different attributes measured on different scales
- Need to be *normalized* (why ?)

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

where  $v_r$  is the actual value of attribute

- Nominal attributes: distance either 0 or 1
- Common policy for missing values: assumed to be 0  
(given normalized attributes)

Add WeChat [edu\\_assist\\_pro.com](https://edu_assist_pro.com)

## When To Consider Nearest Neighbour

# Assignment Project Exam Help

- Inst

- Les

<https://eduassistpro.github.io>

- Lots of training data

- No requirement for “explanatory” model to be le

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

## When To Consider Nearest Neighbour

# Assignment Project Exam Help

Advantages:

- Sta
  - Can <https://eduassistpro.github.io>
  - Training is very fast
  - Can learn complex target functions
  - Don't lose information by generalization - keep
- Add WeChat edu\_assist\_pro

## When To Consider Nearest Neighbour

# Assignment Project Exam Help

Disadvantages:

- Slow at query time: basic algorithm scans entire training data to derive the class label
- "Curse of dimensionality": irrelevant attributes
  - Remedy: attribute selection or weights
- Assumption of linear decision boundaries
- Problem of noisy instances:
  - Remedy: remove from data set
  - not easy – how to know which are noisy ?

Add WeChat `edu_assist_pro`

## When To Consider Nearest Neighbour

# Assignment Project Exam Help

- What is the inductive bias of  $k$ -NN?
- an assumption that the classification of query instance  $x_q$  will be most accurate

$k$ -NN uses

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

- Regression approximating a real-valued target variable
- Residual the error  $\hat{y}(x) - f(x)$  in approximating  $y$
- Kernel function function of distance used to weight training example, i.e., kernel function is the function  $K$  s.t.  
 $w_i = K(d(x_i, x_q))$

Add WeChat `edu_assist_pro`

## Nearest-neighbour classifier

• kNN uses the training data as exemplars, so training is  $O(n)$ , but prediction is also  $O(n)!$ )

- 1N
- By in dec <https://eduassistpro.github.io>
- Easily adapted to real-valued targets, and even to structured objects (nearest-neighbour retrieval). Can also output  $k > 1$
- Warning: in high-dimensional spaces everything and so pairwise distances are uninformative (curse of dimensionality)

Add WeChat edu\_assist\_pro

Distance-Weighted *k*NN

# Assignment Project Exam Help

- Might want to weight nearest neighbours more heavily ...

- Use distance function to construct a weight  $w_i$
- Rep

<https://eduassistpro.github.io>

$$\hat{f}(x_q) \leftarrow \arg \max_{v \in V} w_i \delta(v, f(x))$$

Add WeChat edu\_assist\_pro

where

$$w_i \equiv \frac{1}{d(x_q, x_i)^2}$$

and  $d(x_q, x_i)$  is distance between  $x_q$  and  $x_i$

Distance-Weighted  $k$ NN

# Assignment Project Exam Help

For real-valued target functions replace the final line of the algorithm by:

<https://eduassistpro.github.io>

(denominator normalizes contribution of individual)

Now we can consider using *all* the training examples

- using all examples (i.e., when  $k = n$ ) via  $w_i = \frac{y_i}{n}$
- Shepard's method

Add WeChat  $edu\_assist\_pro$

## Evaluation

# Assignment Project Exam Help

Lazy learners do not construct an explicit model so how do we evaluate the output of the learning process ?

- 1-N
- $k$ -N
- <https://eduassistpro.github.io>

Leave-one-out cross-validation (LOOCV) – leave one out  
predict it given the rest

Add WeChat edu\_assist\_pro

$$(x_1, y_1), (x_2, y_2), \dots, (x_{i-1}, y_{i-1}), \quad i+1 \quad i+1 \quad n \quad n$$

Error is mean over all predicted examples. Fast – no models to be built !

# Curse of Dimensionality

Bellman (1960) coined this term in the context of dynamic programming

Imagine instances described by 20 attributes, but only 2 are relevant to target function — “similar” examples will appear “distant”.

*Curse of dimensionality*

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

One approach:

- Stretch  $j$ th axis by weight  $z_j$  where prediction error
- Use cross-validation to automatically choose weights  $z_1, \dots, z_n$
- Note setting  $z_j$  to zero eliminates this dimension altogether

## Curse of Dimensionality

# Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat `edu_assist_pro`

- number of “cells” in the instance space grows exponentially in the number of features
- with exponentially many cells we would need exponentially many data points to ensure that each cell is sufficiently populated to make nearest-neighbour predictions reliably

## Curse of Dimensionality

# Assignment Project Exam Help

Some ideas to address this for instance-based ( $k$ -nearest-neighbour) learning

- Euclidean distance with weights on attributes

---

<https://eduassistpro.github.io>

- updating of weights based on nearest neighbour

- class correct/incorrect: weight increased

- can be useful if not all features used in classification

:Add WeChat `edu_assist_pro`

See Moore and Lee (1994) "Efficient Algorithms for Minimizing Cross Validation Error"

## Instance-based (nearest-neighbour) learning

# Assignment Project Exam Help

Recap – Practical problems of 1-NN scheme:

- Slow (but fast  $k = d$  tree-based approaches exist)
- Noisy data
- All attributes deemed equally important
  - Remedy: attribute weighting (or simply scale)
- Doesn't perform explicit generalization
  - Remedy: rule-based or tree-based NN approaches

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

Refinements of instance-based (nearest-neighbour) classifiers \*

# Assignment Project Exam Help

- Edi
- ma
- Sav
- IB2: *incremental* NN learner: only incorporates misclassified instances into the classifier
  - Problem: noisy data gets incorporated
- IB3: store *classification performance* i.e.  
& only use in prediction if above a threshold

Add WeChat **edu\_assist\_pro**

## Dealing with noise \*

# Assignment Project Exam Help

use larger values of  $k$  (why ? How to find the "right"  $k$  ?)

- One way: cross-validation-based  $k$ -NN classifier (but slow)
- Diff

kee  
(IB

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

- Computes confidence interval for an instance's success rate and for default accuracy of its class
- If lower limit of first interval is above upper limit of success rate, instance is accepted (IB3: 5%-level)
- If upper limit of first interval is below lower limit of success rate, instance is rejected (IB3: 12.5%-level)

Add WeChat edu\_assist\_pro

## Uncertainty

# Assignment Project Exam Help

As  
cert

<https://eduassistpro.github.io>

–Albert Einstein

Add WeChat edu\_assist\_pro

## Two Roles for Bayesian Methods

# Assignment Project Exam Help

Provides practical learning algorithms:

- Naive Bayes classifier learning
- Bay
- Co
- Ho

<https://eduassistpro.github.io>

Provides useful conceptual framework:

- # Add WeChat edu\_assist\_pro
- Provides a “gold standard” for evaluating other methods
  - Some additional insight into Occam’s razor

## Bayes Theorem

# Assignment Project Exam Help

---

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

$P(h)$  = prior probability of hypothesis  $h$

$P(D)$  = prior probability of training data

$P(h|D)$  = probability of  $h$  given  $D$

$P(D|h)$  = probability of  $D$  given  $h$

Add WeChat edu\_assist\_pro

## Choosing Hypotheses

# Assignment Project Exam Help

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

General

Maximu

<https://eduassistpro.github.io>

$$h_{MAP} = \arg \max_{h \in H}$$

Add WeChat edu\_assist\_pr

$$= \arg \max_{h \in H}$$

$$= \arg \max_{h \in H} P(D|h)P(h)$$

## Choosing Hypotheses

# Assignment Project Exam Help

If assume

Maximu

<https://eduassistpro.github.io>

$$h_{ML} = \arg \max_{h_i \in H} P(D | h_i)$$

Add WeChat edu\_assist\_pro

# Applying Bayes Theorem

Does patient have cancer or not?

A patient takes a lab test and the result comes back positive.

The test returns a correct positive result in only 98% of the cases  
in which

result is positive  
can

Add WeChat edu\_assist\_pro

$$P(\text{cancer}) =$$

$$P($$

$$P(\oplus \mid \text{cancer}) =$$

$$P(\ominus \mid \text{cancer}) =$$

$$P(\oplus \mid \neg \text{cancer}) =$$

$$P(\ominus \mid \neg \text{cancer}) =$$

## Applying Bayes Theorem

Does patient have cancer or not?

# Assignment Project Exam Help

A patient takes a lab test and the result comes back positive.

The test returns a correct positive result in only 98% of the cases

in wh

resu

pre

can

## Add WeChat edu\_assist\_pr

$$P(\text{cancer}) = .008$$

$$P($$

$$P(\oplus | \text{cancer}) = .98$$

$$P(\ominus | \text{cancer}) = .02$$

$$P(\oplus | \neg\text{cancer}) = .03$$

$$P(\ominus | \neg\text{cancer}) = .97$$

## Applying Bayes Theorem

# Assignment Project Exam Help

Does patient have cancer or not?

We can fin

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

$$P(\oplus \mid \neg \text{cancer})P(\neg \text{cancer}) = 0.03 \cdot 0.992 = 0.02976$$

Add WeChat edu\_assist\_pro

Thus  $h_{MAP} = \dots$

## Applying Bayes Theorem

# Assignment Project Exam Help

Does patient have cancer or not?

We can find the maximum a posteriori (MAP) hypothesis

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

Thus  $h_{MAP} = \neg \text{cancer}$ .

Also note: posterior probability of hypothesis *cancer* higher than prior.

## Applying Bayes Theorem

# Assignment Project Exam Help

How to get the posterior probability of a hypothesis  $h$ :

Divide by  $P(\oplus)$ , probability of data, to normalize result for  $h$ :

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

Denominator ensures we obtain posterior probabil

Sum for all possible numerator values, since hypotheses exclusive (e.g., patient either has cancer or does not).

Marginal likelihood (marginalizing out over hypothesis) = probability of the data.

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

## Basic Formulas for Probabilities

*Product Rule:* probability  $P(A \wedge B)$  of conjunction of two events A and B:

**Assignment Project Exam Help**

$$P(A \wedge B) = P(A|B)P(B) = P(B|A)P(A)$$

*Sum Rule*

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

*Theorem of total probability:* if events  $A_1$

with  $\sum_{i=1}^n P(A_i) = 1$ , then:

**Add WeChat edu\_assist\_pro**

$$P(B) = \sum_{i=1}^n P(B|A_i)P(A_i)$$

## Basic Formulas for Probabilities

# Assignment Project Exam Help

Also worth remembering

- *Conditional Probability:* probability of  $A$  given  $B$ :

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

- Rearrange sum rule to get:

Add WeChat edu\_assist\_pro

Exercise: Derive Bayes Theorem.

## Brute Force MAP Hypothesis Learner

# Assignment Project Exam Help

- For e

<https://eduassistpro.github.io>

- Output the hypothesis  $h_{MAP}$  with the line

Add WeChat `edu_assist_pro`

$$h_{MAP} = \arg \max_{h \in H} P(h|D)$$

## Relation to Concept Learning (i.e., classification)

# Assignment Project Exam Help

Canonical concept learning task:

- inst
- con
- fro

$D$   
<https://eduassistpro.github.io>

"zero-error" classification rules)

What would Bayes rule produce as the MAP hypothesis?

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

Does this algorithm output a MAP hypothesis??

## Relation to Concept Learning

# Assignment Project Exam Help

Brute Force MAP Framework for Concept Learning:  
Assume fixed set of instances  $x_1, \dots, x_m$

Assume

Choose

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

- $P(h) = \frac{1}{|H|}$  for all  $h$  in  $H$

Choose  $P(D|h)$ :

- $P(D|h) = 1$  if  $h$  consistent with  $D$
- $P(D|h) = 0$  otherwise

## Relation to Concept Learning

# Assignment Project Exam Help

Then:

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

$$P(h|D) =$$

Add WeChat <sup>l 0 other</sup> edu\_assist\_pro

## Relation to Concept Learning

Note that since likelihood is zero if  $h$  is inconsistent then the posterior is also zero. But how did we obtain the posterior for consistent  $h$ ?

<https://eduassistpro.github.io>

Add WeChat  $\overline{\frac{1}{|V|}}$   
edu\_assist\_pro

$$= \overline{|VS_{H,D}|}$$

## Relation to Concept Learning

How did we obtain  $P(D)$ ? From theorem of total probability:

# Assignment Project Exam Help

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

$$h_i \in VS_{H,D} \quad /$$

Add WeChat  $\sum_{h_i \in VS_{H,D}} \frac{1}{|H|}$   $edu\_assist\_pr$

$$= \frac{|VS_{H,D}|}{|H|}$$

## Evolution of Posterior Probabilities

# Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat `edu_assist_pro`

(a)

hypotheses

(b)

hypotheses

c

## Relation to Concept Learning

# Assignment Project Exam Help

Every hypothesis consistent with  $D$  is a MAP hypothesis, if we assume

- unif
- targ
- det
- etc. (see above)

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

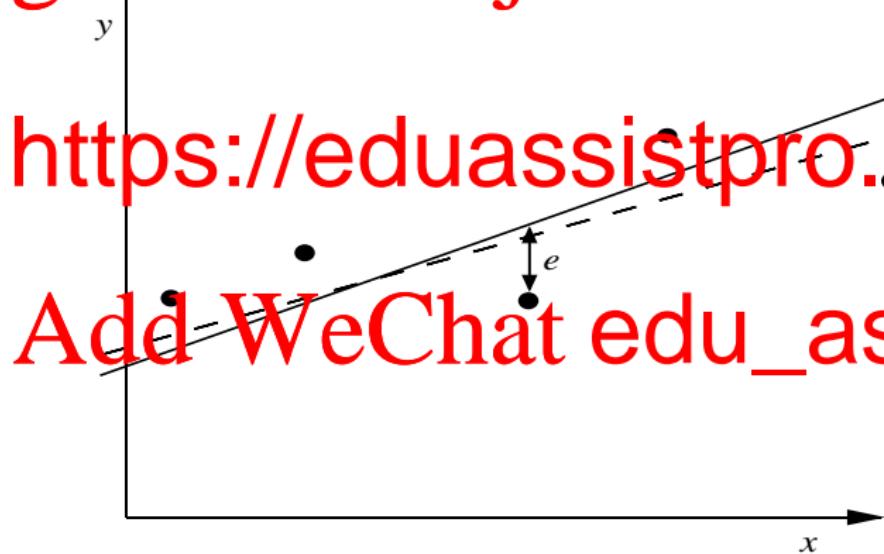
So, this learning algorithm will output a MA  
does not explicitly use *probabilities* in learni

Add WeChat `edu_assist_pro`

# Learning A Real Valued Function

E.g., learning a linear target function  $f$  from noisy examples:

Assignment Project Exam Help



<https://eduassistpro.github.io>  
Add WeChat edu\_assist\_pro

## Learning A Real Valued Function

# Assignment Project Exam Help

Consider any real-valued target function  $f$

Training Examples  $\langle x_i, d_i \rangle$ , where  $d_i$  is noisy training value

- $d_i$
- $e_i$  i  
acc

<https://eduassistpro.github.io>

Then the **maximum likelihood** hypothesis  
minimizes the sum of squared errors:

Add WeChat  $_{\text{m}}$  [edu\\_assist\\_pro](https://edu_assist_pro)

$$h_{ML} = \arg \min_{h \in H} \sum_{i=1}^m (d_i - h(x_i))^2$$

## Learning A Real Valued Function

How did we obtain this ?

# Assignment Project Exam Help

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

$$h \in H \quad i=1$$

$$\arg\max_{h \in H} \prod_{i=1}^m \frac{1}{\sqrt{2}}$$

Add WeChat edu\_assist\_pro

Recall that we treat each probability  $p(D|h)$  as if  $h = f$ , i.e., we assume  $\mu = f(x_i) = h(x_i)$ , which is the key idea behind maximum likelihood !

## Learning A Real Valued Function

Maximize natural log to give simpler expression:

$$h_{ML} = \arg \max_{h \in H} \ln \frac{1}{\sqrt{\pi^2}} - \frac{1}{2} \frac{(d_i - h(x_i))^2}{\sigma^2}$$

<https://eduassistpro.github.io>

$$h_{ML} = \arg \max_{h \in H} \sum_{i=1}^m -(d_i - h(x_i))^2$$

Equivalently, we can minimize the positive version of the expression:

$$h_{ML} = \arg \min_{h \in H} \sum_{i=1}^m (d_i - h(x_i))^2$$

# Discriminative and generative probabilistic models

- *Discriminative models* model the posterior probability distribution  $P(Y|X)$ , where  $Y$  is the target variable and  $X$  are the features. That is, given  $X$  they return a probability distribution over  $Y$ .

# Assignment Project Exam Help

- *Generative models* model the joint distribution  $P(Y, X)$  of the target  $Y$  and the fe

deri

In par

erior

distribution can be obtained as  $P(Y|X) = \frac{P(Y, X)}{P(X)}$ .

- Alternatively, generative models can be described as  $P(X|Y)$ , since  $P(Y, X) = P(X|Y)P(Y)$ . The prior ( $P(X)$ , usually abbreviated to ‘prior’) can be easily estimated.
- Such models are called ‘generative’ because we can sample from the joint distribution to obtain new data points together with their labels. Alternatively, we can use  $P(Y)$  to sample a class and  $P(X|Y)$  to sample an instance for that class.

# Assessing uncertainty in estimates

Suppose we want to estimate the probability  $\theta$  that an arbitrary e-mail is spam, so that we can use the appropriate prior distribution.

## Assignment Project Exam Help

- The natural thing to do is to inspect  $n$  e-mails, determine the number of spam e-mails  $d$ , and set  $\hat{\theta} = d/n$ ; we don't really need any com

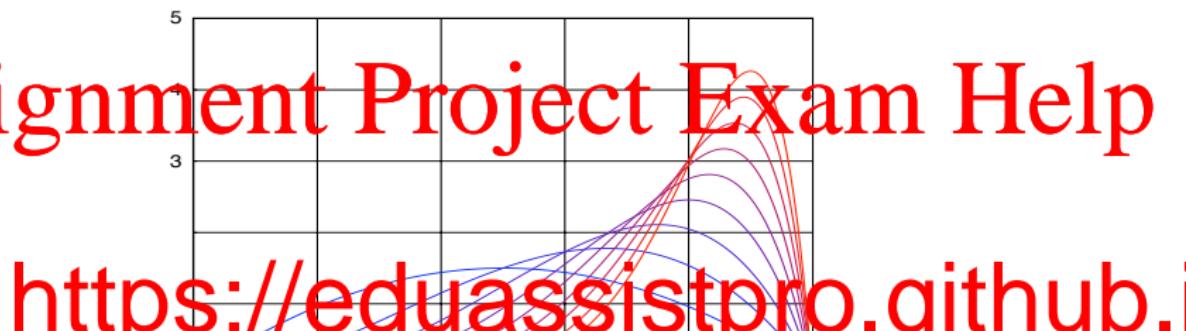
- Ho https://eduassistpro.github.io/ximuma pos are completely ruled out.

- We model this by a probability distribution over  $\theta$  (prior distribution in this case) which is updated each time new information

This is further illustrated in the figure for a distribution that is more and more skewed towards spam.

- For each curve, its bias towards spam is given by the area under the curve and to the right of  $\theta = 1/2$ .

# Assessing uncertainty in estimates



Each time we inspect an e-mail, we are reducing our uncertainty about the prior spam probability  $\theta$ . After we inspect two e-mails and both turn out to be spam, the possible  $\theta$  values are characterised by a symmetric distribution around 1/2. If we inspect a third, fourth, ..., tenth e-mail and each time (except the first one) it is spam, then this distribution narrows and shifts a little bit to the right each time. The distribution for  $n$  e-mails reaches its maximum at  $\hat{\theta}_{\text{MAP}} = \frac{n-1}{n}$  (e.g.,  $\hat{\theta}_{\text{MAP}} = 0.8$  for  $n = 5$ ).

# The Bayesian perspective

Assignment Project Exam Help

Explicitly modelling the posterior distribution over the parameter  $\theta$  has a number of advantages that are usually associated with the "Bayesian" perspective.

- We can estimate <https://eduassistpro.github.io>.
- We can obtain a generative model for the parameter by sampling from the posterior distribution, which contains more information than a summary statistic such as the mean – so, rather than using a single e-mail with a summary statistic, a generative model can contain a number of e-mails with  $\theta$  sampled from the posterior distribution.

## The Bayesian perspective

- We can quantify the probability of statements such as 'e-mails are biased toward ham' (the tiny shaded area in the figure demonstrates that after observing one ham and nine spam e-mails this probability is very low)
- We can also quantify the probability of statements such as 'e-mails are symmetric' (the large shaded area in the figure on the previous slide) if we know our prior.

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

one in the figure on the previous slide) as our prior.

The key point is that probabilities do not have to be interpreted as estimates of relative frequencies, but can carry the more subjective (possibly subjective) degrees of belief.

Consequently, we can attach a probability distribution to almost anything: not just features and targets, but also model parameters and even models.

# Minimum Description Length Principle

# Assignment Project Exam Help

Once again, the MAP hypothesis

$$h_{MAP} = \arg \max P(D|h)P(h)$$

Which is eq

<https://eduassistpro.github.io>

$$h_{MAP} = \arg \max_{h \in H} \log_2 P($$

Or

Add WeChat edu\_assist\_pro

$$h_{MAP} = \arg \min_{h \in H} -\log_2 P(D|h) - \log_2 P(h)$$

## Minimum Description Length Principle

# Assignment Project Exam Help

Interestingly, this is an expression about a quantity of *bits*.

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

From information theory:

The optimal (shortest expected coding length) with probability  $p$  is  $-\log_2 p$  bits.

Add WeChat `edu_assist_pr`

## Minimum Description Length Principle

# Assignment Project Exam Help

So interpr

: - lo  
: - lo <https://eduassistpro.github.io>

**Note well:** assumes *optimal* encodings, w

are known. In practice, this is difficult, and makes a differ

Add WeChat edu\_assist\_pro

## Minimum Description Length Principle

# Assignment Project Exam Help

Occam's razor: prefer the shortest hypothesis

MDL: pre-

<https://eduassistpro.github.io>

$$MDL_{h \in H}$$

$$C_1 \quad C_2$$

where  $L_{C_i}(\cdot)$  is the description length of  $x$

Add WeChat edu\_assist\_pro

## Minimum Description Length Principle

# Assignment Project Exam Help

Without loss of generality, classifier is here assumed to be a decision tree.

Example:  $H = \text{decision trees}$ ,  $D = \text{training data labels}$

- $L_{C_1}$
- $L_{C_2}$
- $C_2$  describe exceptions
- Hence  $h_{MDL}$  trades off tree size for training error  
i.e., prefer the hypothesis that minimizes

$$\text{length}(h) + \text{length}(\text{misclassifications})$$

## Most Probable Classification of New Instances

# Assignment Project Exam Help

So far we've

$h_{MAP}$ )

$D$  (i.e.,

<https://eduassistpro.github.io>

Given new instance  $x$ , what is its most probable *classification*?

- $h_{MAP}(x)$  is not the most probable classification

Add WeChat edu\_assist\_pro

## Most Probable Classification of New Instances

# Assignment Project Exam Help

Consider

- Thr
- Giv

$$h_1(x) = +, h_2(x) = -, h_3(x) = -$$

- What's most probable classification of

<https://eduassistpro.github.io>

## Bayes Optimal Classifier

# Assignment Project Exam Help

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

Add WeChat `edu_assist_pro`

$P(h_1|D) = .4, P(-|h_1) = 0$   
 $P(h_2|D) = .3, P(-|h_2) = 1$   
 $P(h_3|D) = .3, P(-|h_3) = 1, P(+|h_3) = 0$

## Bayes Optimal Classifier

therefore  
**Assignment Project Exam Help**

$$P(+|h_i)P(h_i|D) = .4$$

<https://eduassistpro.github.io>

and

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

*No other classification method using the same hypothesis space and same prior knowledge can outperform this method on average*

# Bayes Error

What is the best performance attainable by a (two-class) classifier ?

Defining the probability of error for classifying some instance  $x$  by

# Assignment Project Exam Help

$$P(\text{error } x) = P(\text{class}_1 | x) \quad \text{if we predict class}_2$$

<https://eduassistpro.github.io>

This gives

$$\sum_x P(\text{error}) = P(\text{er})$$

So we can justify the use of the decision rule

if  $P(\text{class}_1 | x) > P(\text{class}_2 | x)$  then predict class<sub>1</sub>  
else predict class<sub>2</sub>

*On average, this decision rule minimises probability of classification error.*

## Naive Bayes Classifier

# Assignment Project Exam Help

When to use

- Most effective classifier
- Attributes that describe instances are conditionally independent given classification

Successful applications:

Add WeChat edu\_assist\_pro

- Classifying text documents
- Gaussian Naive Bayes for real-valued data

## Naive Bayes Classifier

Assume target function  $f: X \rightarrow V$ , where each instance is described by attributes  $\langle a_1, a_2 \dots a_n \rangle$ .

Most prob

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

$$\begin{aligned} v_{MAP} &\equiv \arg \max_{v_i \in V} \frac{P(a_1, a_i)}{P} \\ &= \arg \max_{v_j \in V} P(a_1, a_j) \end{aligned}$$

Add WeChat edu\_assist\_pro

## Naive Bayes Classifier

Naive Bayes assumption:

# Assignment Project Exam Help

$$P(a_1, a_2 \dots a_n | v_j) = P(a_i | v_j)$$

- Att <https://eduassistpro.github.io>
  - which means knowledge about the value of a particular attribute tells us nothing about the value of another attribute

which gives [Add WeChat edu\\_assist\\_pro](#)

**Naive Bayes classifier:**  $v_{NB} = \arg \max_{v_j \in V} P(v_j) \prod_i P(a_i | v_j)$

# Naive Bayes Algorithm

Naive\_Bayes\_Learn(*examples*)

For each target value  $v_j$

$P(v_j) \leftarrow \text{estimate } P(v_j)$

# Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pro

Classify\_New\_Instance( $x$ )

$$v_{NB} = \arg \max_{v_j \in V} \hat{P}(v_j) \prod_{a_i \in x} \hat{P}(a_i | v_j)$$

## Naive Bayes Example

# Assignment Project Exam Help

Consider PlayTennis again. .

									dy	es
Sun	2	3								
Ove	4	0								
Rai	3	2								
Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5	False	6/9
Overcast	4/9	0/5	Mild	4/9	2/5	Nor	6	1/		2/5
Rainy	3/9	2/5	Cool	3/9	1/5					
Play										
Yes	9	5								
No										
	9/14	5/14								

Add WeChat edu\_assist\_pro

## Naive Bayes Example

# Assignment Project Exam Help

Say we have the new instance:

We want to  
(<https://eduassistpro.github.io>)  
true)

$v_{NB} = \arg \max_{v_j \in \{ "yes", "no" \}} P$   
Add WeChat edu\_assist\_pr

## Naive Bayes Example

So we first calculate the likelihood of the two classes, "yes" and "no"

# Assignment Project Exam Help

for "yes"

$$P(\text{true}|y)$$

0.00

<https://eduassistpro.github.io>

for "no"

$$= P(n) \times P(\text{sun}|n) \times P(\text{no true}|n)$$
$$= \frac{5}{14} \times \frac{3}{5} \times \frac{1}{5} \times \frac{4}{5} \times \frac{3}{5}$$
$$= 0.0206$$

## Naive Bayes Example

# Assignment Project Exam Help

Then convert to a probability by normalisation

<https://eduassistpro.github.io>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

The Naive Bayes classification is “no”.

## Naive Bayes: Subtleties

Conditional independence assumption is often violated

# Assignment Project Exam Help

- ...b

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

$$\arg \max_{v_j \in V} \hat{P}(v_j) \quad \hat{P}(a_i | v_j) = \arg \max_{v_j \in V} )$$

i.e. maximum probability is assigned to correct class

Add WeChat edu\_assist\_pro

- see [Domingos & Pazzani, 1996] for analysis
- Naive Bayes posteriors often unrealistically close to 1 or 0
- adding too many redundant attributes will cause problems (e.g. identical attributes)

Naive Bayes: “zero-frequency” problem

# Assignment Project Exam Help

What if none of the training instances with target value  $v$  have attribute value  $a_i$

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

$$P(v_j) \quad P(a_i|v_j) = 0$$

Pseudo-counts (add 1 to each count) (version of the Add WeChat  $\overset{i}{\text{edu\_assist\_pr}}$ )

Naive Bayes: “zero-frequency” problem

# Assignment Project Exam Help

- In some cases adding a constant different from 1 might be more appropriate

- Exa

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

$$\frac{\bar{2}}{9+\mu} \quad \frac{\bar{3}}{9+\mu} \quad \frac{\bar{3}}{9+\mu}$$

- Weights don't need to be equal (if they sum to 1) – a for prior

Add WeChat edu\_assist\_pro

$$\frac{2+\mu p_1}{9+\mu} \quad \frac{4+\mu p_2}{9+\mu} \quad \frac{3+\mu p_3}{9+\mu}$$

Naive Bayes: “zero-frequency” problem

# Assignment Project Exam Help

This generalisation is a Bayesian estimate for  $\hat{P}(v_i|v_j)$

$$\frac{n_c + mp}{n}$$

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

- $n$  is number of training examples for which  $v = v_j$ ,
- $n_c$  number of examples for which  $v = v_i$
- $p$  is prior estimate for  $\hat{P}(v_i|v_j)$
- $m$  is weight given to prior (i.e. number of “virtual” examples)

This is called the  $m$ -estimate of probability.

Add WeChat `edu_assist_pro`

## Naive Bayes: missing values

# Assignment Project Exam Help

- Tra  
valu <https://eduassistpro.github.io>
- Classification: attribute will be omitted from ca

Add WeChat edu\_assist\_pro

## Naive Bayes: numeric attributes

Assignment Project Exam Help

Usual assumption: attributes have a normal or Gaussian probability distribution (given the class)

- The mean  $\mu$  is given by the average of the values  $x_i$ :

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

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

The standard deviation  $\sigma$ :

$$\sigma = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2}$$

Naive Bayes: numeric attributes

# Assignment Project Exam Help

Then we have the density function  $f(x)$ :

$$\frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2}}$$

Example <https://eduassistpro.github.io>  
standard deviation = 6.2. Density value

$$f(\text{temperature} = 60 | \text{"yes"}) = \frac{1}{\sqrt{2\pi} \cdot 6.2}$$

Missing values during training are not included in calculation of mean and standard deviation.

## Naive Bayes: numeric attributes

Note: the normal distribution is based on the simple exponential function

# Assignment Project Exam Help

$$f(x) = e^{-|x|^m}$$

As the power function.

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

Where  $m = 2$

$$f(x) = e^{-|x|}$$

and this is the basis of the normal distribution – the variance is the result of scaling so that the integral (the area under the  $-\infty$  to  $+\infty$ ) is equal to 1.

from “Statistical Computing” by Michael J. Crawley (2002) Wiley.

## Categorical random variables

Categorical variables or features (also called discrete or nominal) are ubiquitous in machine learning.

Assignment Project Exam Help

- Perhaps the most common form of the Bernoulli distribution models whether or not a word occurs in a document. That is, for the  $i$ -th word  $X_i$ , the probability of occurrence is denoted by a Bernoulli parameter  $\theta_i$ .

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

- Variables with more than two outcomes are also known as categorical variables. For example, every word position in an e-mail corresponds to a categorical variable with  $c$  outcomes, where  $c$  is the vocabulary size. The multinomial distribution manifests itself as a histogram of the number of occurrences of all vocabulary words in a document. This establishes an alternative way of modelling text documents that allows the number of occurrences of a word to influence the classification of a document.

Add WeChat edu\_assist\_pro

## Categorical random variables

Both these document models are in common use. Despite their differences, they both assume independence between word occurrences, generally referred to as the *naive Bayes assumption*.

- In the multinomial document model, this follows from the very use of the word distribution vector.
- In the Bernoulli document model, the words in the document vector are statistically independent, which allows us to calculate the joint probability of a particular bit vector as the product of the probabilities of each component.
- In practice, such word independence assumptions are often not true: if we know that an e-mail contains the word 'Viagra', we can be quite sure that it will also contain the word 'pill'. Violated independence assumptions reduce the quality of probability estimates but may still allow good classification performance.

## Example application: Learning to Classify Text

# Assignment Project Exam Help

In machine  
learning t

<https://eduassistpro.github.io>

Here is a simplified version in the multinomial document

Add WeChat edu\_assist\_pro

## Learning to Classify Text

# Assignment Project Exam Help

Why?

- Lea
- Lea

<https://eduassistpro.github.io>

Naive Bayes is among most effective algorithms

Add WeChat edu\_assist\_pro

What attributes shall we use to represent text documents?

## Learning to Classify Text

# Assignment Project Exam Help

Target concept *Interesting?* : Document  $\rightarrow \{+, -\}$

① Rep

② Lea

- $P(+)$

- $P(-)$

- $P(doc|+)$

- $P(doc|-)$

• <https://eduassistpro.github.io>

• Add WeChat edu\_assist\_pro

## Learning to Classify Text

# Assignment Project Exam Help

Naive Bayes conditional independence assumption

<https://eduassistpro.github.io>

where  $P(a_i = w_k | v_j)$  is probability that word in p

$v_j$

one more assumption:  $P(a_i = w_k | v_j) = P(w_k)$

Add WeChat edu\_assist\_pro

"bag of words"

## Learning to Classify Text

LEARN\_NAIVE\_BAYES\_TEXT(*Examples*, *V*)

// collect all words and other tokens that occur in Examples

# Assignment Project Exam Help

*Vocabulary*  $\leftarrow$  all distinct words and other tokens in *Examples*

// calculate

for each topic

*docs<sub>j</sub>*

$$P(v_j) \leftarrow \frac{|docs_j|}{|Examples|}$$

*Text<sub>j</sub>*  $\leftarrow$  a single document created by concatenating

*n*  $\leftarrow$  total number of words in *Text<sub>j</sub>* (count

for each word *w<sub>k</sub>* in *Vocabulary*

*n<sub>k</sub>*  $\leftarrow$  number of times word *w<sub>k</sub>* occurs in *Text<sub>j</sub>*

$$P(w_k|v_j) \leftarrow \frac{n_k + 1}{n + |Vocabulary|}$$

## Learning to Classify Text

# Assignment Project Exam Help

CLASSIFY\_NAIVE\_BAYES\_TEXT(*Doc*)

- pos
- Voc
- Ret

nd in  
<https://eduassistpro.github.io>

$v_{NB} = \arg \max_{v_i \in V} P(v_i)$

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

## Application: 20 Newsgroups

Given: 1000 training documents from each group

Learning task: classify each new document by newsgroup it came from

Assignment Project Exam Help

comp.graphics

misc.forsale

<https://eduassistpro.github.io>

alt.atheism

sci.space

soc.religion.christian

sci.crypt

talk.religion.misc

sci.electro

talk.politics.mideast

sci.med

talk.politics.misc

talk.politics.guns

Naive Bayes: 89% classification accuracy

## Article from rec.sport.hockey

Path: cantaloupe.srv.cs.cmu.edu!das-news.harvard.edu!ogies!uwm.edu

From: xxxyyyzzz.edu (John Doe)

Subject: Re: This year's biggest and worst (opinion)...

Date: 5 Apr 9

I can only hope  
for pleasure

a defensive defenseman, but he's clearly much more than that.

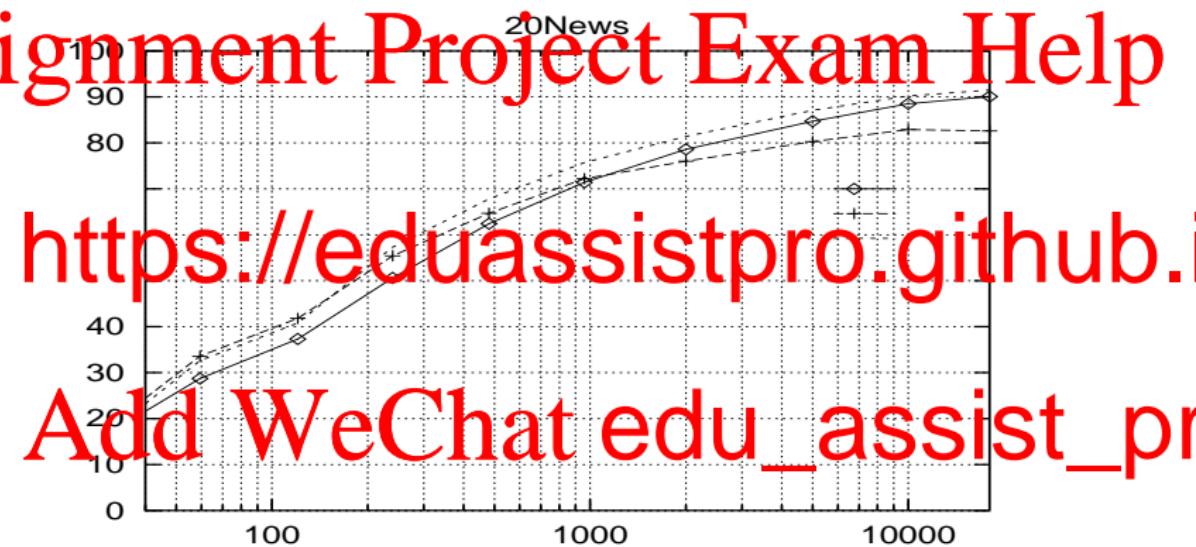
Great skater and hard shot (though wish he were more accurate)

In fact, he pretty much allowed the Kings to trade away their

huge defensive liability Paul Corley. Kelly Hrudey

biggest disappointment if you thought he was any good to begin with. But, at best, he's only a mediocre goaltender. A better choice would be Tomas Sandstrom, though not through any fault of his own, but because some thugs in Toronto decided ...

## Learning Curve for 20 Newsgroups



Accuracy vs. Training set size (1/3 withheld for test)

## Probabilistic decision rules

We have chosen one of the possible distributions to model our data  $X$  as coming from either class

# Assignment Project Exam Help

- The more different  $P(X | Y = \text{spam})$  and  $P(X | Y = \text{ham})$  are, the more effective the decision rules:
- Thus, the maximum likelihood (ML) – predict  $\arg \max_{Y \in \{\text{spam}, \text{ham}\}} P(Y)$

maximum likelihood (ML) – predict  $\arg \max_{Y \in \{\text{spam}, \text{ham}\}} P(Y)$

maximum a posteriori (MAP) – predict

$$\arg \max_y P(X = x | Y = y) P(Y = y);$$

The relation between the first two decision rules is that ML classification is equivalent to MAP classification with a uniform class distribution.

Probabilistic decision rules

# Assignment Project Exam Help

We again  
illustrate

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

Add WeChat edu\_assist\_pro

## Prediction using a naive Bayes model

Suppose our vocabulary contains three words  $a$ ,  $b$  and  $c$ , and we use a multivariate Bernoulli model for our e-mails, with parameters

)

This means  
(+), com

<https://eduassistpro.github.io>

The e-mail to be classified contains words

described by the bit vector  $\mathbf{x} = (1, 1, 0)$ . We

Add WeChat edu\_assist\_pr

$$P(\mathbf{x}|\oplus) = 0.5 \cdot 0.67 \cdot$$

$$P(\mathbf{x}|\ominus) = 0.67 \cdot 0.33 \cdot (1 - 0.33) = 0.148$$

The ML classification of  $\mathbf{x}$  is thus spam.

## Prediction using a naive Bayes model

In the case of two classes it is often convenient to work with likelihood ratios and odds.

# Assignment Project Exam Help

- The likelihood ratio can be calculated as

$$\frac{P(\mathbf{x} | \oplus)}{P(\mathbf{x} | \ominus)} = \frac{0.5}{0.67} \approx 0.74$$

- This odds are more than 2/3, but ham if they are less than that.
- For example, with 33% spam and 67% ham, resulting in a posterior odds of  $\frac{P(\oplus)}{P(\ominus)} = \frac{0.33}{0.67} = 1/2$ , resulting in a posterior odds of 1/2.

Add WeChat edu\_assist\_pro

$$\frac{P(\oplus | \mathbf{x})}{P(\ominus | \mathbf{x})} = \frac{P(\mathbf{x} | \oplus) P(\oplus)}{P(\mathbf{x} | \ominus) P(\ominus)} = 3/2 \cdot 1/2 = 3/4 < 1$$

In this case the likelihood ratio for  $\mathbf{x}$  is not strong enough to push the decision away from the prior.

## Prediction using a naive Bayes model

Alternatively, we can employ a multinomial model. The parameters of a multinomial establish a distribution over the words in the vocabulary, say

# Assignment Project Exam Help

The e-mail

occurred

by the count

occurred

a, one single

described

<https://eduassistpro.github.io>

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

$$P(\mathbf{x}|\oplus) = 4! \frac{0.3^3}{3!} \frac{0.5^1}{1!} \frac{0.2^0}{0!}$$

$$P(\mathbf{x}|\ominus) = 4! \frac{0.6^3}{3!} \frac{0.2^1}{1!} \frac{0.2^0}{0!}$$

The likelihood ratio is  $(\frac{0.3}{0.6})^3 (\frac{0.5}{0.2})^1 (\frac{0.2}{0.2})^0 = 5/16$ . The ML classification of  $\mathbf{x}$  is thus ham, the opposite of the multivariate Bernoulli model. This is mainly because of the three occurrences of word  $a$ , which provide strong evidence for ham.

## Training data for naive Bayes

A small e-mail data set described by count vectors.

# Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat `edu_assist_pro`

	3			
$e_4$	2	3	0	
$e_5$	4	3	0	
$e_6$	4	0	3	
$e_7$	3	0	0	—
$e_8$	0	0	0	—

## Training data for naive Bayes

The same data set described by bit vectors

# Assignment Project Exam Help

---

---

<https://eduassistpro.github.io>

	3			
$e_4$	1	1	0	
$e_5$	1	1	0	
$e_6$	1	0	1	
$e_7$	1	0	0	—
$e_8$	0	0	0	—

## Training a naive Bayes model

Consider the following e-mails consisting of five words  $a, b, c, d, e$ :

$e_1: d \; d \; e \; b \; b \; d \; e$

$e_5: a \; b \; a \; b \; a \; b \; a \; e \; d$

$e_2$

$e_3$

$e_4$

<https://eduassistpro.github.io>

We are told that the e-mails on the left are spam and those on the right are ham, and so we use them as a small training set to train our classifier.

Add WeChat edu\_assist\_pr

- First, we decide that  $d$  and  $e$  are so-called common to convey class information.
- The remaining words,  $a, b$  and  $c$ , constitute our vocabulary.

## Training a naive Bayes model

For the multinomial model, we represent each e-mail as a count vector, as before.

# Assignment Project Exam Help

- In order to train the classifier, we need to estimate the parameters  $\theta_{\oplus}$  and  $\theta_{\ominus}$  for the spam and ham classes respectively. We can do this by counting the frequency of each word in the training set. For example, if we have 11 words in the vocabulary, we might end up with the following counts:  
<https://eduassistpro.github.io/nb.html>
- To start, we will consider the case where each e-mail has exactly one word. This means that each e-mail is represented by a single word, which brings the total number of words in the vocabulary down to 20 for each class.
- The estimated parameter vectors are thus:  
 $\hat{\theta}_{\oplus} = (6/20, 10/20, 4/20) = (0.3, 0.5, 0.2)$  for spam and  
 $\hat{\theta}_{\ominus} = (12/20, 4/20, 4/20) = (0.6, 0.2, 0.2)$  for ham.

Add WeChat edu\_assist\_pro

## Training a naive Bayes model

In the multivariate Bernoulli model e-mails are represented by bit vectors, as before.

# Assignment Project Exam Help

- Adding the bit vectors for each class results in  $(2, 3, 1)$  for spam and  $(3,$
- Each word particular vocabulary word.
- Probability smoothing now means adding two containing each word and one containing none of
- This results in the estimated parameter vectors $\hat{\theta}^{\oplus} = (3/6, 4/6, 2/6) = (0.5, 0.67, 0.33)$  for spam and $\hat{\theta}^{\ominus} = (4/6, 2/6, 2/6) = (0.67, 0.33, 0.33)$  for ham.

# Linear discriminants

# Assignment Project Exam Help

Many forms of linear discriminant from statistics and machine learning,  
e.g.,

- Fish
- Log
- Perceptron
  - a linear threshold classifier
  - an early version of an artificial “neuron”
  - still a useful method, and source of ideas

# Logistic regression

In the case of a two-class problem, model the probability of one class

**Assignment Project Exam Help**

1

or

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

$$\ln \frac{P(Y = 1|\mathbf{x})}{1 - P(Y = 1|\mathbf{x})}$$

The quantity on the l.h.s. is called the logit model for the logit.

Note: to fit this is actually more complex than linear regression, so we omit the details.

Generalises to multiple class versions ( $Y$  can have more than two values).

# Perceptron

A linear classifier that can achieve perfect separation on linearly separable data is the *perceptron*, originally proposed as a simple *neural network* by T. Rosenblatt in the late 1950s.

# Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pro

## Perceptron

Originally implemented in software (based on the McCulloch-Pitts neuron from the 1940s), then in hardware as a 20x20 visual sensor array with potentiometers for adaptive weights.

# Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

Source <http://en.wikipedia.org/w/index.php?curid=47541432>

# Assignment Project Exam Help

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

Add WeChat `edu_assist_pro`

Output  $o$  is thresholded sum of products of inputs and their weights:

$$o(x_1, \dots, x_n) = \begin{cases} +1 & \text{if } w_0 + w_1x_1 + \dots + w_nx_n > 0 \\ -1 & \text{otherwise.} \end{cases}$$

# Assignment Project Exam Help

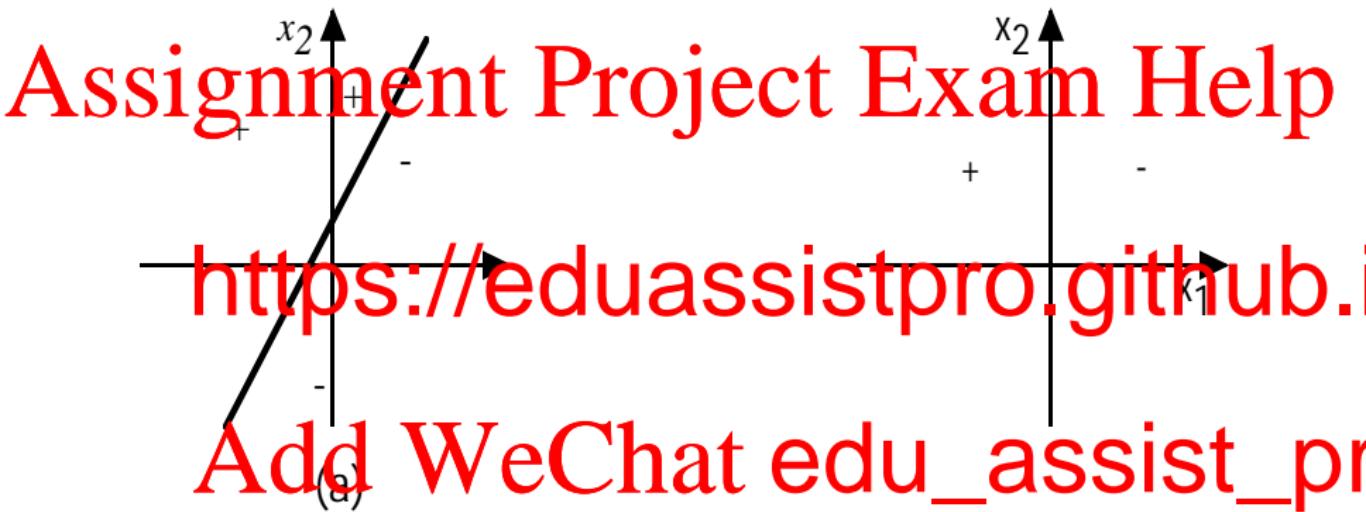
<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

Or in vector notation:

$$o(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w} \cdot \mathbf{x} > 0 \\ -1 & \text{otherwise.} \end{cases}$$

## Decision Surface of a Perceptron



Represents some useful functions

- What weights represent  $o(x_1, x_2) = AND(x_1, x_2)$ ?
- What weights represent  $o(x_1, x_2) = XOR(x_1, x_2)$ ?

## Decision Surface of a Perceptron

# Assignment Project Exam Help

So some functions not representable

- e.g.

• <https://eduassistpro.github.io/>

- for non-linearly separable data we'll need some

• e.g., networks of these ...

• the start of “deep” networks ...

• Add WeChat edu\_assist\_pro

# Perceptron learning

Key idea:

# Assignment Project Exam Help

Perceptr

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

where the weight update  $\Delta w_i$  depends only o

les and

is modulated by a “smoothing” parameter

“learning rate”.

Add WeChat edu\_assist\_pro

Can prove that perceptron learning will converge:

- if training data is linearly separable
- and  $\eta$  sufficiently small

# Perceptron learning

The perceptron iterates over the training set, updating the weight vector every time it encounters an incorrectly classified example.

# Assignment Project Exam Help

- For example, let  $\mathbf{x}_i$  be a misclassified positive example, then we have

$y_i = 1$  and  $\mathbf{w} \cdot \mathbf{x}_i < t$ , which means that

$$\mathbf{w}' = \mathbf{w} + \eta \mathbf{x}_i$$

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

- This can be achieved by calculating the new weight vector

$\mathbf{w}' = \mathbf{w} + \eta \mathbf{x}_i$ , where  $0 < \eta \leq 1$  is the learning rate. We then have  $\mathbf{w}' \cdot \mathbf{x}_i = \mathbf{w} \cdot \mathbf{x}_i + \eta \mathbf{x}_i \cdot \mathbf{x}_i$ . Since  $\mathbf{x}_i \cdot \mathbf{x}_i > 0$ , we have  $\mathbf{w}' \cdot \mathbf{x}_i > \mathbf{w} \cdot \mathbf{x}_i$ .

- Similarly, if  $\mathbf{x}_j$  is a misclassified negative example, then we have  $y_j = -1$  and  $\mathbf{w} \cdot \mathbf{x}_j > t$ . In this case we calculate the new weight vector as  $\mathbf{w}' = \mathbf{w} - \eta \mathbf{x}_j$ , and thus  $\mathbf{w}' \cdot \mathbf{x}_j = \mathbf{w} \cdot \mathbf{x}_j - \eta \mathbf{x}_j \cdot \mathbf{x}_j < \mathbf{w} \cdot \mathbf{x}_j$ .

## Perceptron learning

- The two cases can be combined in a single update rule:

$$\mathbf{w}' = \mathbf{w} + \eta y_i \mathbf{x}_i$$

- Her  
when <https://eduassistpro.github.io>
- Thi  
classification
- The algorithm just iterates over the training set a  
weight update rule until all the examples are corre
- If there is a linear model that separates the positive from the negative  
examples, i.e., the data is linearly separable, it can be shown that the  
perceptron training algorithm will converge in a finite number of steps.

# Perceptron training algorithm

**Algorithm** Perceptron( $D, \eta$ ) // perceptron training for linear classification

**Input:** labelled training data  $D$  in homogeneous coordinates, learning rate  $\eta$ .

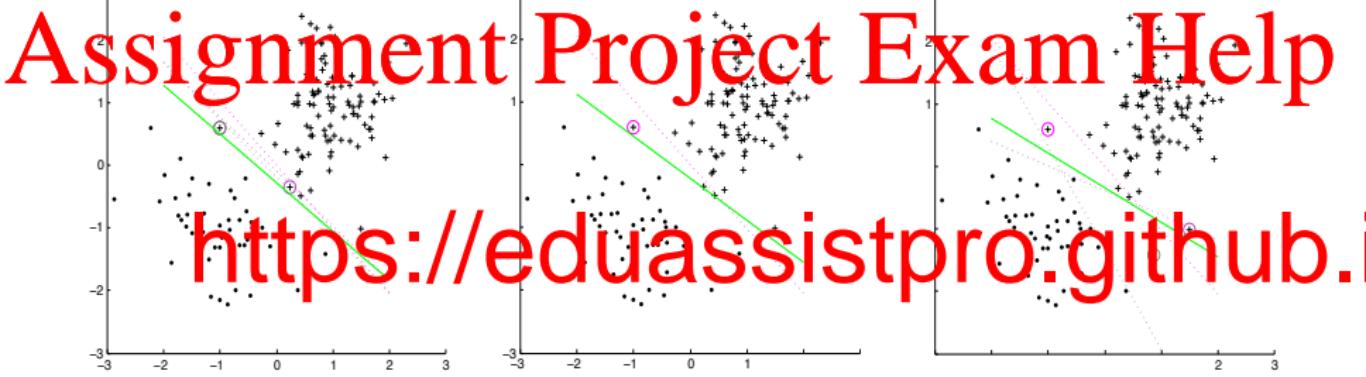
**Output:** weight vector  $w$  defining classifier  $\hat{y} = \text{sign}(w \cdot x)$ .

```
1 w ← 0
2 conver
3 while c
4   con
5   for i = 1 to |D| do
6     if  $y_i w \cdot x_i \leq 0$  then // i.e.,  $\hat{y} \neq y$ 
7       w ← w +  $\eta y_i x_i$ 
8       converged ← false // We changed
9   end
10 end
11 end
```

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pro

## Perceptron training – varying learning rate



Add WeChat `edu_assist_pro`

(left) A perceptron trained with a small learning rate (circled examples are the ones that trigger the weight update. (middle) Increasing the learning rate to  $\eta = 0.5$  leads in this case to a rapid convergence. (right) Increasing the learning rate further to  $\eta = 1$  may lead to too aggressive weight updating, which harms convergence. The starting point in all three cases was the basic linear classifier.

## Summary

# Assignment Project Exam Help

- Two major frameworks for classification by linear models

Dist

Pro

- We have seen that the perceptron is a form of threshold model.
- These classifiers are also, in some sense, linear models.
- So we have established the basis for learning classification models.
- Later we will see how to extend by building on these ideas

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)