

Assignment Project Exam Help

Introduction to Machine Learning and Data Mining

<https://eduassistpro.github.io>

Add WeChat [edu_assist_pro](#)

Last revision: 28 Feb 2021

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/~tom/mlbook.html>

Material derived from slides for the book

"Bayesian Reasoning and Machine Learning" by D. Barber

Cambridge University Press (2012)

<http://www.cs.toronto.edu/~dbarber/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, giving an overview of some of the key ideas and approaches we will cover in the course.

Following it you should be able to describe some of the main concepts and outline so

- cat
- wid
- batch vs. online settings
- parametric vs. non-parametric approaches
- generalisation in machine learning
- training, validation and testing phases in applications
- limits on learning

Add WeChat `edu_assist_pro`

What we will cover

Assignment Project Exam Help

- cor
- fou
- rele
- practical applications

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

What we will NOT cover

Assignment Project Exam Help

- lots o
- lots o
- “bi <https://eduassistpro.github.io>
- commercial and business aspects of “analytics”
- ethical aspects of AI and ML

Add WeChat edu_assist_pro
although all of these are interesting and important topics

Some history

Assignment Project Exam Help

One can imagine that after the machine had been in operation for some time, the instructions would have been altered out of recognition, but nevertheless still be such that one would have to ad

*calc
desi*

<https://eduassistpro.github.io>

*efficient manner. In such a case one would have to ad
the progress of the machine had not been foresee
original instructions were put in. It would be like a p
learnt much from his master, but had added much
own work.*

Add WeChat edu_assist_pro

From A. M. Turing's lecture to the London Mathematical Society. (1947)

Some definitions

Assignment Project Exam Help

The field of machine learning is concerned with the question of how to construct computer programs that automatically improve from experience.

"Machine L

<https://eduassistpro.github.io>

Machine learning, then, is about making computers modify or adapt their actions (whether these actions are making predictions or controlling a robot) so that these actions are more accurate, where accuracy is measured by how often the chosen actions reflect the correct ones.

"Machine Learning". S. Marsland (2015)

Some definitions

Machine learning is the systematic study of algorithms and systems that improve their knowledge or performance with experience.

Machine Le

The term “machine learning” refers to the ability of computers to automatically learn and identify meaningful patterns in data.

“Understanding Machine Learning” S. Shalev-Shwartz and

Add WeChat edu_assist_pro

Data mining is the extraction of implicit, previously unknown, and potentially useful information from data.

“Data Mining”. I. Witten et al. (2016)

Machine Learning is . . .

Assignment Project Exam Help

Trying to g

<https://eduassistpro.github.io>

R. Kohn (2015)

Add WeChat edu_assist_pro

How is Machine Learning different from . . .

Machine learning comes originally from Artificial Intelligence (AI), where the motivation is to build intelligent agents, capable of acting autonomously. Learning is a characteristic of intelligence, so to be successful, one must understand and control the agent's environment.

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

These are not requirements in:

- statistics — the results are typically mathematical models
- data mining — the results are typically models of the world for humans

These criteria are often also necessary, but not always sufficient, for machine learning.

Machine Learning for Human-level Artificial Intelligence

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

ML for human-level AI, right ?

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

ML for human-level AI, right ? Not so fast . . .

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Supervised and unsupervised learning

Assignment Project Exam Help

The most widely used categories of machine learning algorithms are:

- *Sup*
- *Uns*

<https://eduassistpro.github.io>

There are a

strategies to acquire data, such as reinforcement learning.

Add WeChat edu_assist_pro

Note: output class can be real-valued or discrete, scalar structure ...

Supervised and unsupervised learning

Assignment Project Exam Help

Supervised learning tends to dominate in applications.

Why ?

<https://eduassistpro.github.io>

Generally, because it is much easier to define the problem and develop an error measure (loss function) to evaluate different alg settings, data transformations, etc. for supervised i.e. unsupervised learning.

Add WeChat edu_assist_pro

Supervised and unsupervised learning

Assignment Project Exam Help

Unfortunately ...

In the real w

sufficien

<https://eduassistpro.github.io>

So in such cases unsupervised learning is really what you

Add WeChat edu_assist_pro

but currently, finding good unsupervised learning or
machine learning tasks remains a research challenge

Machine learning models

Assignment Project Exam Help

Machine learning models can be distinguished according to their main intuition

- *Geometric* (hy) <https://eduassistpro.github.io>
- *Probabilistic* models view learning as a process of reducing uncertainty, modelled by means of probability
- *Logical* models are defined in terms of easily interpretable expressions.

Add WeChat edu_assist_pro

Machine learning models

Assignment Project Exam Help

Alternatively, can be characterised by *algorithmic properties*:

- *Reg*
- *Cla*
- *Neu*
- *Local models* predict in the local region of a query
- *Tree-based models* partition the data to make predictions
- *Ensembles* learn multiple models and combine them

Add WeChat **edu_assist_pro**

Linear regression

Given 2 real-valued variables X_1, X_2 , labelled with a real-valued variable Y , find “line of best fit” that captures the dependency of Y on X_1, X_2 .

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Learning here is by minimizing MSE, i.e., the average of the squared vertical distances of values of Y from the learned function $\hat{Y} = \hat{f}(\mathbf{X})$.

Linear regression

A question: is it possible to do better than the line of best fit?

Assignment Project Exam Help

Maybe. Linear regression assume that the (x_i, y_i) examples in the data are "generat

So any trial $y = f(x)$

But what if f is non-linear ?

Add WeChat edu_assist_pro

We may be able to reduce the mean squared error (MSE) v
 $\sum_i (y_i - \hat{y})^2$ by trying a different function.

Can “decompose” MSE to aid in selecting a better function.

Nearest Neighbour

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Nearest Neighbour is a regression or classification algorithm that predicts whatever is the output value of the nearest data point to some query.

Classification

Customer103: (time=t0)

Years of credit: 9
Loan balance: \$2,400
Income: \$52k
Own Ho

Customer103: (time=t1)

Years of credit: 9
Loan balance: \$3,250
Income: ?

...

Customer103: (time=tn)

Years of credit: 9
Loan balance: \$4,500
Income: ?

<https://eduassistpro.github.io>

If $OtherDelinquentAccounts > 2$, and

$NumberDelinquentBillingCycles > 1$

Then $ProfitableCustomer = No$

If $OtherDelinquentAccounts = 0$, and

$Income > 30k$ OR $YearsOfCredit > 3$

Then $ProfitableCustomer = Yes$

Assassinating spam e-mail

SpamAssassin is a widely used open-source spam filter. It calculates a score for an incoming e-mail, based on a number of built-in rules or ‘tests’ in SpamAssassin’s terminology, and adds a junk flag and a summary report to the e-mail’s headers if the score is 5 or more.

-0.1 RCVD_IN_M

| | |
|------------------------|--|
| 0.6 HTML_IMA | |
| 1.2 TVD_FW_GI | BODY: HTML in |
| 0.0 HTML_MESS | BODY: HTML font face is not a word |
| 0.6 HTML_FONX_FACE_BAD | FULL: Email has an inline gif |
| 1.4 SARE_GIF_ATTACH | MTA bounce message |
| 0.1 BOUNCE_MESSAGE | Message is some kind of bounce message |
| 0.1 ANY_BOUNCE_MESSAGE | AWL from: address is not in auto-whitelist |
| 1.4 AWL | |

Add WeChat edu_assist_pr

From left to right you see the score attached to a particular test, the test identifier, and a short description including a reference to the relevant part of the e-mail. As you see, scores for individual tests can be negative (indicating evidence suggesting the e-mail is ham rather than spam) as well as positive. The overall score of 5.3 suggests the e-mail might be spam.

Linear classification

Suppose we have only two tests and four training e-mails, one of which is spam. Both tests succeed for the spam e-mail; for one ham e-mail neither test succeeds, for another the first test succeeds and the second doesn't, and for the third ham e-mail the first test fails and the second succeeds.

It is easy to see
these four e

introduced above we could describe this classifier as $4x_1 + 4x_2 > 5$ or
 $(4, 4) \cdot (x_1, x_2) > 5$.

Add WeChat edu_assist_pro

In fact, any weight between 2.5 and 5 will ensure that the threshold is only exceeded when both tests succeed. We could even consider assigning different weights to the tests – as long as each weight is less than 5 and their sum exceeds 5 – although it is hard to see how this could be justified by the training data.

Spam filtering as a classification task

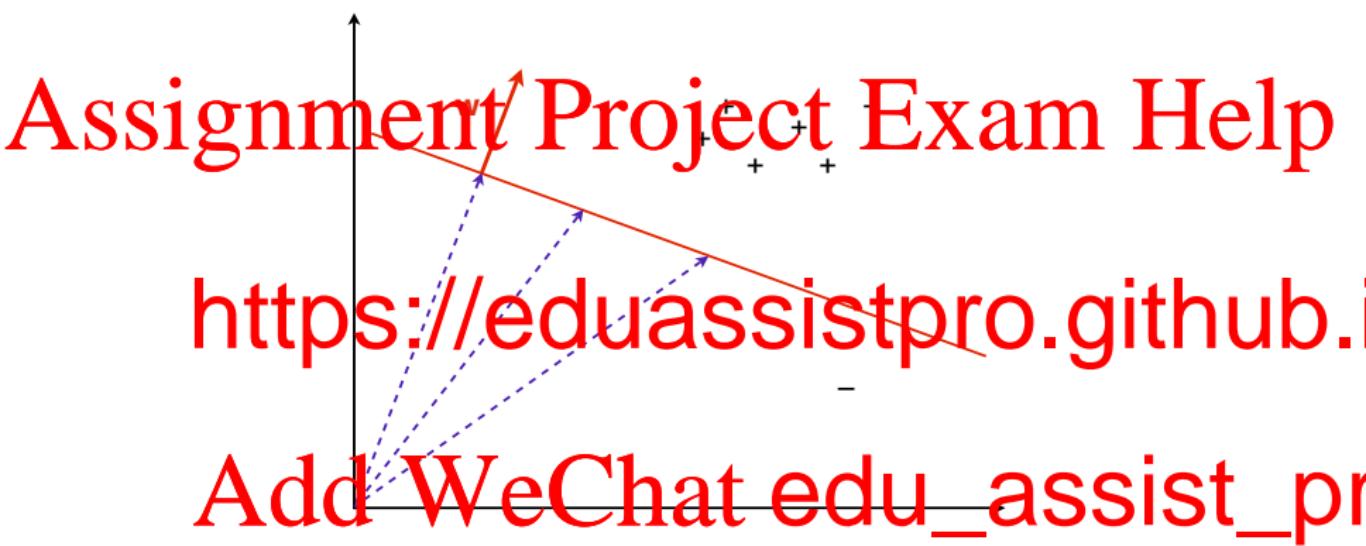
The columns marked x_1 and x_2 indicate the results of two tests on four different e-mails. The fourth column indicates which of the e-mails are spam. The right-most column demonstrates that by thresholding the function

<https://eduassistpro.github.io>

| E-mail | x_1 | x_2 | Spam? | $4x$ | $+ 4x$ |
|--------|-------|-------|-------|------|--------|
| 1 | 1 | 1 | 1 | 8 | |
| 2 | 0 | 0 | 0 | 0 | |
| 3 | 1 | 0 | 0 | 4 | |
| 4 | 0 | 1 | 0 | 4 | |

Add WeChat edu_assist_pro

Linear classification in two dimensions



- straight line separates positives from negatives
- defined by $\mathbf{w} \cdot \mathbf{x}_i = t$
- \mathbf{w} is perpendicular to decision boundary
- \mathbf{w} points in direction of positives
- t is the decision threshold

Linear classification in two dimensions

Assignment Project Exam Help

Note: x_i points in the $\mathbf{w} \cdot \mathbf{x}_0 = \|\mathbf{w}\| \|\mathbf{x}_0\| = t$ (where $\|\mathbf{x}\|$ denotes the length of the vector \mathbf{x}).

Add WeChat edu_assist_pro

Homogeneous coordinates

It is sometimes convenient to simplify notation further by introducing an extra constant variable $x_0 = 1$, the weight of which is fixed to $w_0 = -t$.

The exten

ended

weight ve

$w^\circ \cdot x^\circ$

<https://eduassistpro.github.io>

Thanks to these so-called *homogeneous coo*

dary

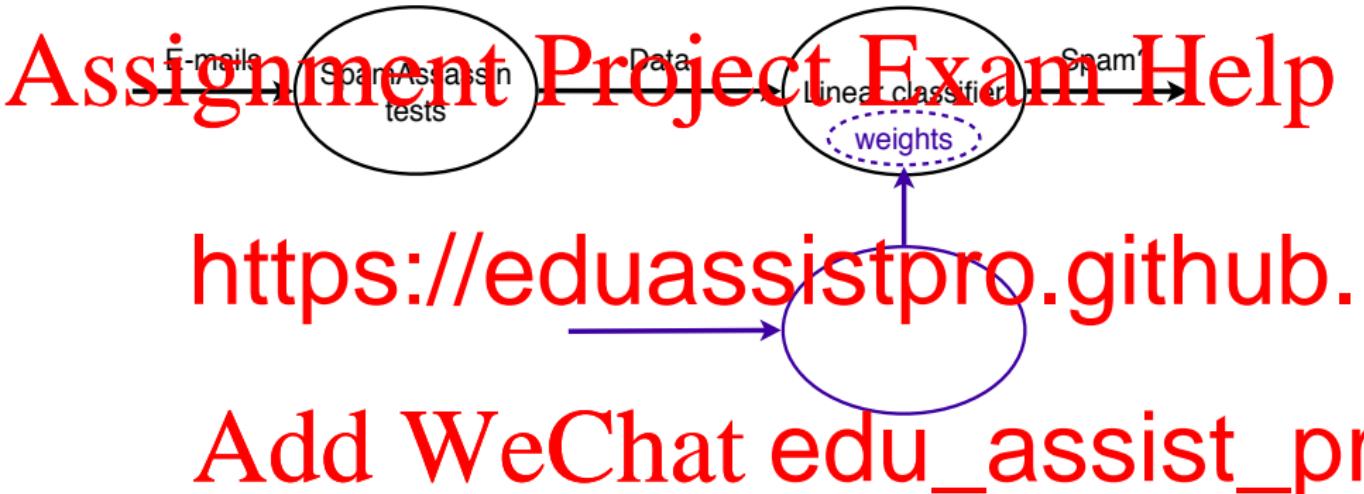
passes through the origin of the extended coordinate s

expense of needing an additional dimension.

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

Note: this doesn't really affect the data, as all data points and the 'real' decision boundary live in the plane $x_0 = 1$.

Machine learning for spam filtering



At the top we see how SpamAssassin approaches the spam e-mail classification task: the text of each e-mail is converted into a data point by means of SpamAssassin's built-in tests, and a *linear classifier* is applied to obtain a 'spam or ham' decision. At the bottom (in blue) we see the bit that is done by machine learning.

A Bayesian classifier I

Bayesian spam filters maintain a *vocabulary* of words and phrases –

Potential spam pattern indicators – for which statistics are collected from a *training set*.

- For i

e-m
tha

<https://eduassistpro.github.io>

this e-mail is spam are 4:1, or the probability of it being spam is 0.80 and the probability of it being ham is

- The situation is slightly more subtle because we have to take account the prevalence of spam. Suppose that I received one spam e-mail for every six ham e-mails. This means that I would estimate the odds of an unseen e-mail being spam as 1:6, i.e., non-negligible but not very high either.

A Bayesian classifier II

Assignment Project Exam Help

If I then learn that the e-mail contains the word 'Viagra', which

occ

two

sim

pro

<https://eduassistpro.github.io>

In this way you are combining two independent pieces o

concerning the prevalence of spam, and the other conc

occurrence of the word 'Viagra', pulling in opposite dir

Add WeChat edu_assist_pro

A Bayesian classifier III

Assignment Project Exam Help

The nice thing about this ‘Bayesian’ classification scheme is that it can be repeated in favour of s

and suppose the prior odds are 4:1 in favour of spam. Then the combined odds are 4:1 times 3:1 is 12:1, which is ample to overcome the 1:6 odds associated with the low prevalence of spam (to give a spam probability of 0.67, up from 0.10).

A rule-based classifier

Assignment Project Exam Help

If the e-mail contains the word 'Viagra' then estimate the odds of spam as 4:1;

- oth
- of sp
- oth <https://eduassistpro.github.io>

The first rule covers all e-mails containing the word 'Via'

whether they contain the phrase 'blue pill', so no overco

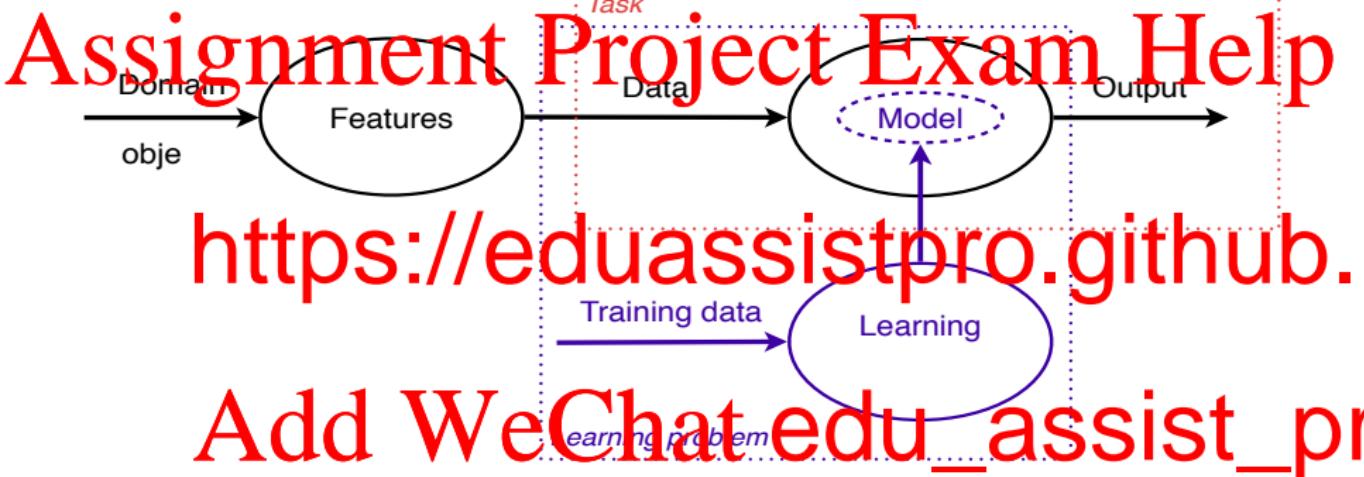
The second rule only covers e-mails containin

not the word 'Viagra', by virtue of the 'otherwise' clause

covers all remaining e-mails: those which neither contain neither 'Viagra' nor 'blue pill'.

Add WeChat `edu_assist_pro`

How machine learning helps to solve a task



An overview of how machine learning is used to address a given task. A task (red box) requires an appropriate mapping – a model – from data described by features to outputs. Obtaining such a mapping from training data is what constitutes a learning problem (blue box).

Some terminology I

Assignment Project Exam Help

Tasks are a
learning al

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Some terminology II

Assignment Project Exam Help

Machine
right mod

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Some terminology III

Assignment Project Exam Help

Models let

give it unity

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Some terminology IV

Assignment Project Exam Help

Does the al

of learnin

gorithm.

<https://eduassistpro.github.io>

If however, it can continue to learn a new data arrives, it is an **online learning** algorithm.

Add WeChat edu_assist_pro

Some terminology V

Assignment Project Exam Help

If the mode

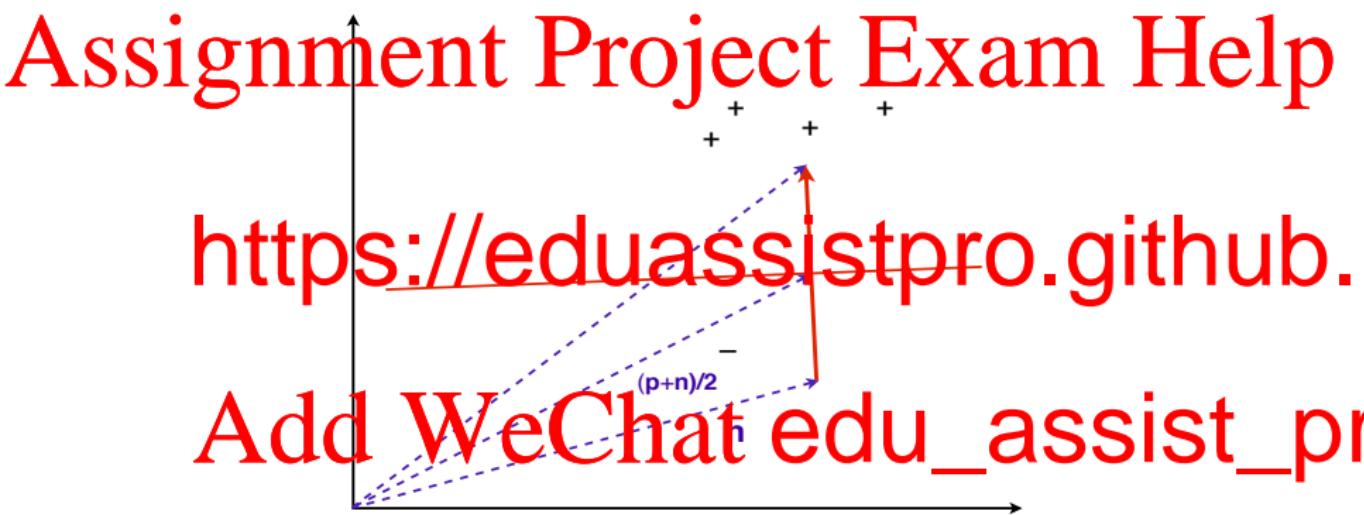
param

<https://eduassistpro.github.io>

Otherwise, if the number of parameters grows with the amount of training data it is categorised as non-parametric.

Add WeChat edu_assist_pro

Basic linear classifier I



The basic linear classifier constructs a decision boundary by half-way intersecting the line between the positive and negative centres of mass.

Basic linear classifier II

Assignment Project Exam Help

The basic li

$w = p$

is on the dec

$t = (p - n) \cdot (p + n)/2 = (||p|| - ||n||)/2$, where $\|x\|$ denotes the length of vector x .

t , with

$(p + n)/2$

Add WeChat edu_assist_pro

Neural Networks I

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Neural Networks II

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Neural Networks III

Assignment Project Exam Help

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

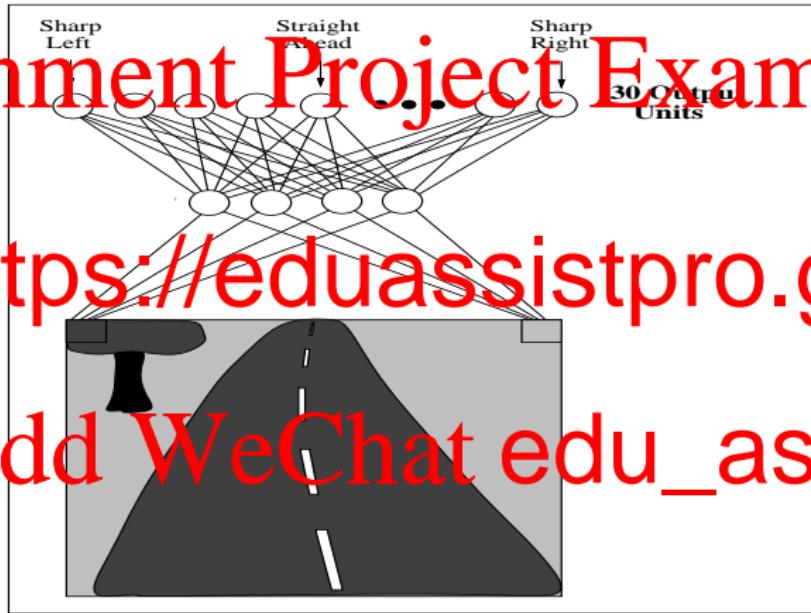
Add WeChat edu_assist_pro

Neural Networks IV

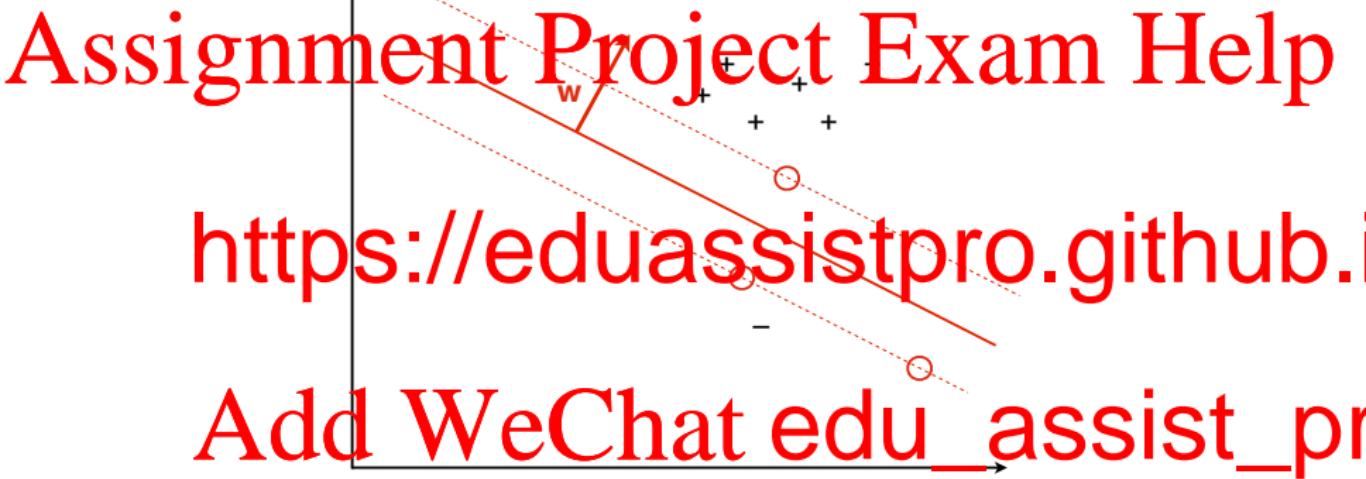
Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



Support vector machine



The decision boundary learned by a support vector machine from the linearly separable data from before. The decision boundary maximises the margin, which is indicated by the dotted lines. The circled data points are the support vectors.

A simple probabilistic model

Assignment Project Exam Help

'Viagra' and 'lottery' are two Boolean features; Y is the class variable, with value

indicate if <https://eduassistpro.github.io> lottery)

| Viagra | lotter | | | |
|--------|--------|------|--|------|
| 0 | 0 | 0.31 | | 0.69 |
| 0 | 1 | 0.65 | | |
| 1 | 0 | 0.80 | | |
| 1 | 1 | 0.40 | | |

Add WeChat edu_assist_pro

Decision rule

Assignment Project Exam Help

Assuming that X and Y are the only variables we know and care about, the posterior distribution $P(Y|X)$ helps us to answer many questions of interest.

- For <https://eduassistpro.github.io>
‘Via probability $P(Y = \text{spam}|\text{Viagra, lotte})$ probability exceeds 0.5 and ham otherwise’
- Such a recipe to predict a value of Y and the posterior distribution $P(Y|X)$ is called a *decision rule*.

Missing values I

Suppose we skimmed an e-mail and noticed that it contains the word 'lottery' but we haven't looked closely enough to determine whether it uses the word 'Viagra'. This means that we don't know whether to use the second or the first definition.

would pre-

and ham if it did

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

The solution is to average these two rows, using the probability occurring in any e-mail (spam or not):

Add WeChat edu_assist_pro

$$\begin{aligned} P(Y|\text{lottery}) &= P(Y|\text{Viagra} = 0, \text{lottery}) \\ &\quad + P(Y|\text{Viagra} = 1, \text{lottery})P(\text{Viagra} = 1) \end{aligned}$$

Missing values II

Assignment Project Exam Help

For instance, suppose for the sake of argument that one in ten e-mails contain th

$P(\text{Viagr}$

$P(Y = s$

$P(Y = \text{ham} | \text{lottery} = 1) = 0.35 \cdot 0.90 + 0$

the

occurrence of 'Viagra' in any e-mail is relatively rare, the r distribution deviates only a little from the second row in

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

Likelihood ratio

As a matter of fact, statisticians work very often with different conditional probabilities, given by the *likelihood function* $P(X|Y)$.

Assignment Project Exam Help

- Like to think of these as thought experiments: If somebody were to send me a spam e-mail, how likely would it be that it contains exactly the word

- **ham**
- **Wh**

<https://eduassistpro.github.io>

- their ratio: how much more likely is it to observe this combination of words in a spam e-mail than it is in a non-spam e-mail

- For instance suppose that for a particular e-mail we have $P(X|Y = \text{spam}) = 3.5 \cdot 10^{-5}$ and $P(X|Y = \text{ham}) = 1.5 \cdot 10^{-6}$,

then observing X in a spam e-mail is nearly five times more likely than it is in a ham e-mail.

- This suggests the following decision rule: predict spam if the likelihood ratio is larger than 1 and ham otherwise.

When to use likelihoods

Assignment Project Exam Help

Use likelihoods
uniform, a

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

Add WeChat edu_assist_pro

Posterior odds

$P(Y = \text{spam}|\text{Viagra} = 0, \text{lottery} = 0) / P(Y = \text{ham}|\text{Viagra} = 0, \text{lottery} = 0)$

<https://eduassistpro.github.io>⁷

$P(Y = \text{ham}|\text{Viagra} = 0, \text{lottery}$

$P(Y = \text{spam}|\text{Viagra} = 1, \text{lotter}$

Add WeChat edu_assist_pro

Using a MAP decision rule we predict ham in the top two cases and spam in the bottom two. Given that the full posterior distribution is all there is to know about the domain in a statistical sense, these predictions are the best we can do: they are *Bayes-optimal*.

Example marginal likelihoods

Assignment Project Exam Help

| Y | $P(\text{Viagra} = 1 Y)$ | $P(\text{Viagra} = 0 Y)$ |
|-----|--------------------------|--------------------------|
| | 0.6 | |
| | 0.88 | |

<https://eduassistpro.github.io>

| Y | $P(\text{lottery} = 1 Y)$ | |
|------|---------------------------|------|
| spam | 0.21 | 0.79 |
| ham | 0.13 | 0.87 |

Add WeChat edu_assist_pro

Using marginal likelihoods

Using the marginal likelihoods from before, we can approximate the likelihood ratios (the previously calculated odds from the full posterior distribution are shown in brackets):

Assignment Project Exam Help

$$\frac{P(\text{Viagra} = 0 | Y = \text{spam})}{P(\text{Viagr})} \frac{P(\text{lottery} = 0 | Y = \text{spam})}{P(\text{Viagr})} \frac{0.60}{0.60} \frac{0.79}{0.79} \frac{.62}{.62} (0.45)$$

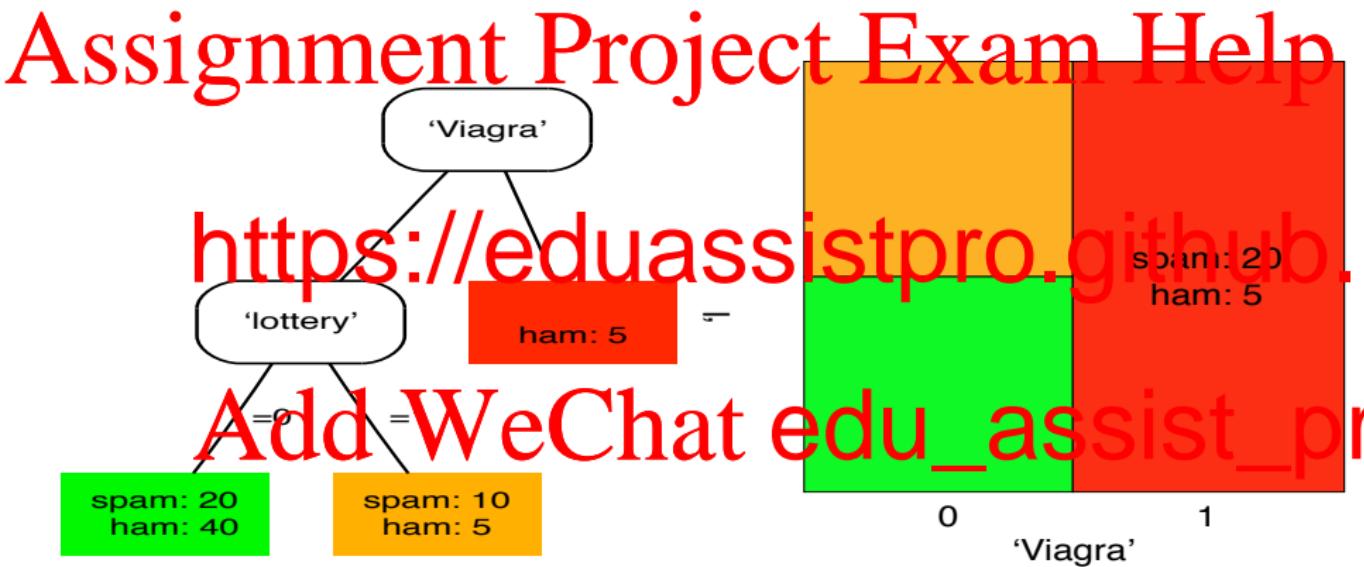
$$\frac{P(\text{Viagr})}{P(\text{Viagr})} \frac{P(\text{lottery} = 0 | Y = \text{spam})}{P(\text{Viagr})} \frac{0.40}{0.40} \frac{0.79}{0.79} (0.45)$$

$$\frac{P(\text{Viagra} = 1 | Y = \text{spam})}{P(\text{Viagra} = 1 | Y = \text{ham})} \frac{P(\text{lottery} = 0 | Y = \text{spam})}{P(\text{lottery} = 0 | Y = \text{ham})} \frac{0.40}{0.40} \frac{0.79}{0.79} (0.45)$$

$$\frac{P(\text{Viagra} = 1 | Y = \text{spam})}{P(\text{Viagra} = 1 | Y = \text{ham})} \frac{P(\text{lottery} = 1 | Y = \text{spam})}{P(\text{lottery} = 1 | Y = \text{ham})} \frac{0.60}{0.60} \frac{0.79}{0.79} (0.45)$$

We see that, using a maximum likelihood decision rule, our very simple model arrives at the *Bayes-optimal* prediction in the first three cases, but not in the fourth ('Viagra' and 'lottery' both present), where the marginal likelihoods are actually very misleading.

A classification tree I



A classification tree combining two Boolean features.

A classification tree II

Assignment Project Exam Help

Each internal node or split is labelled with a feature, and each edge emanating from a split is labelled with a feature value.

Each leaf th

<https://eduassistpro.github.io>

Also indic

training set.

Add WeChat edu_assist_pro

A classification tree partitions the instance space into regions

one for each leaf. We can clearly see that the majority of ham lives in the lower left-hand corner.

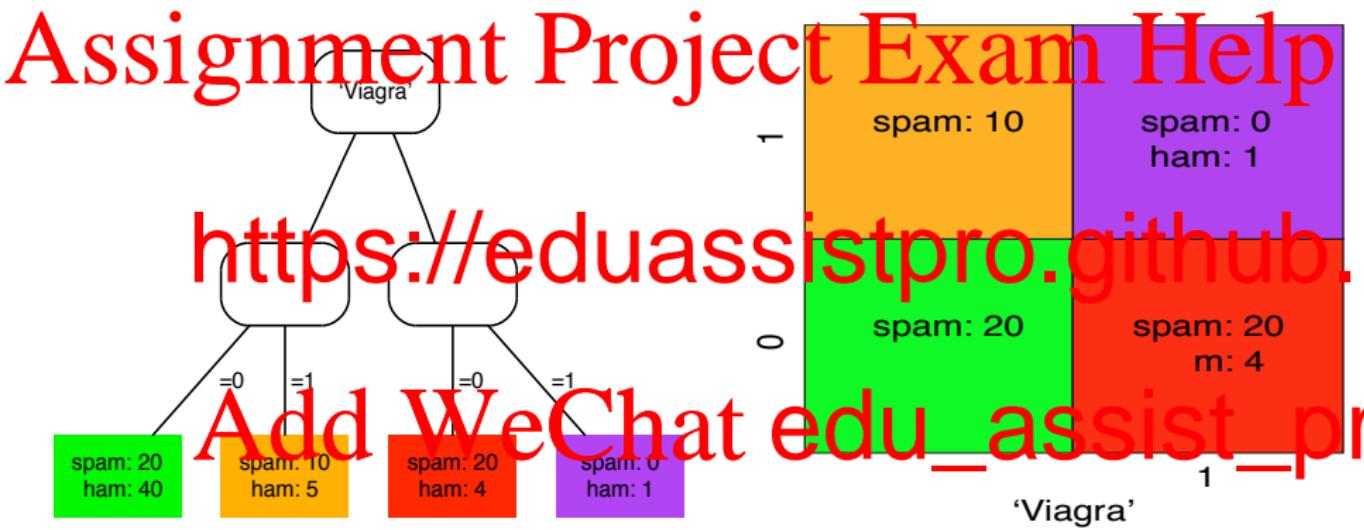
Labelling a classification tree

Assignment Project Exam Help

- The leaves of the classification tree could be labelled, from left to right
maj
- Altering the labels occurring in each leaf: from left to right, $1/3$, $2/3$, and $4/5$.
- Or, if our task was a regression task, we could label the predicted real values or even linear functions of so many real-valued features.

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

A complete classification tree



A complete classification tree built from two Boolean features. The corresponding instance space partition is the finest partition that can be achieved with those two features.

Two uses of features

Assignment Project Exam Help

Suppose

$$1 \leq x \leq 1.$$

A linear ap

$$y = 0.$$

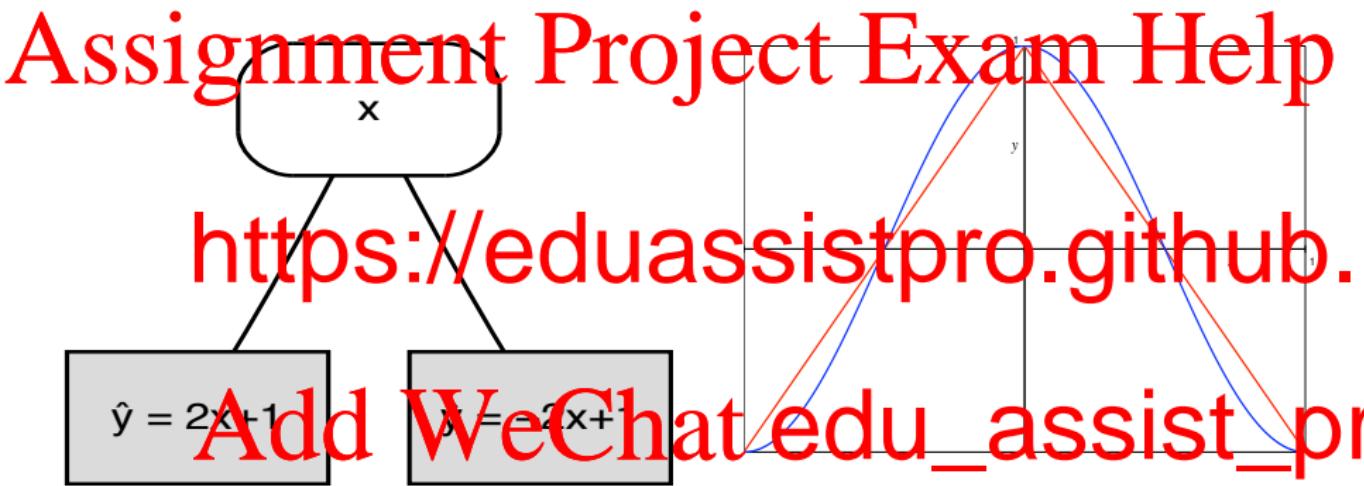
$$0 \leq x \leq$$

interval. We can achieve this by using x b

a regression variable.

Add WeChat edu_assist_pr

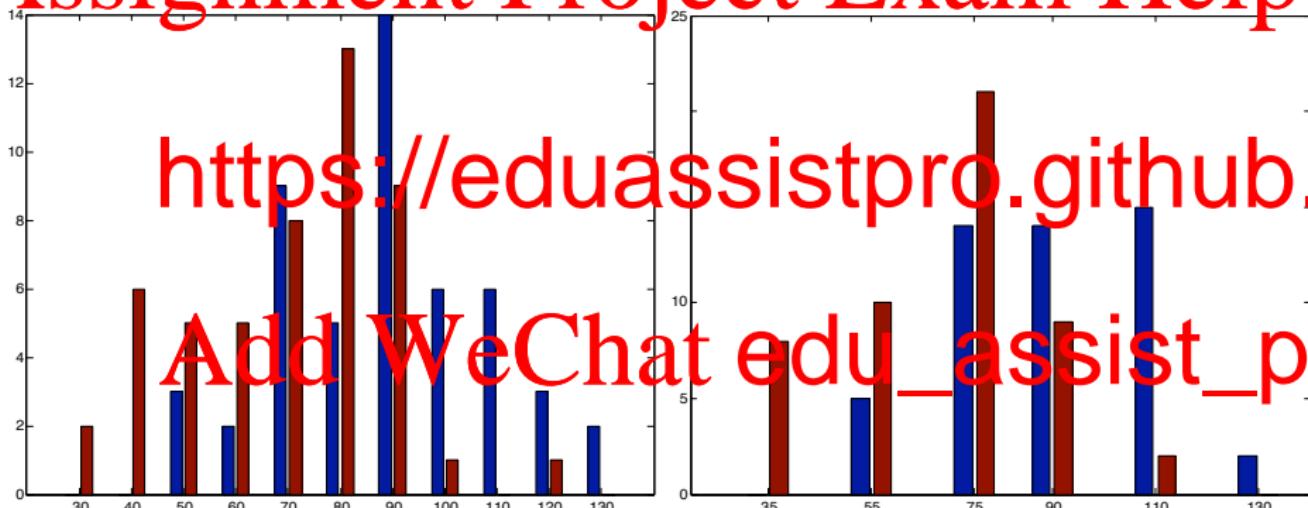
A small regression tree



A regression tree combining a one-split feature tree with linear regression models in the leaves. Note: x used as both a splitting feature and regression variable. At right, function $y = \cos \pi x$ on the interval $-1 \leq x \leq 1$, and piecewise linear approximation by regression tree.

Class-sensitive discretisation I

Assignment Project Exam Help



Class-sensitive discretisation II

Assignment Project Exam Help

Artificial data depicting a histogram of body weight measurements of people with diabetes. The data is split into 10 bins of 10 kilograms each.

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

By joining the

eighth, ninth and tenth intervals, we obtain a discretised feature. The proportion of diabetes cases increases from left to right. This class-sensitive discretisation makes the feature more useful in predicting diabetes.

Add WeChat edu_assist_pro

The kernel trick

Let $\mathbf{x}_1 = (x_1, y_1)$ and $\mathbf{x}_2 = (x_2, y_2)$ be two data points, and consider the mapping $(x, y) \mapsto (x^2, y^2, \sqrt{2}xy)$ to a three-dimensional feature space.

The points in feature space corresponding to \mathbf{x}_1 and \mathbf{x}_2 are

$$\mathbf{x}'_1 = (x_1^2, y_1^2, \sqrt{2}x_1y_1) \quad \text{and} \quad \mathbf{x}'_2 = (x_2^2, y_2^2, \sqrt{2}x_2y_2)$$

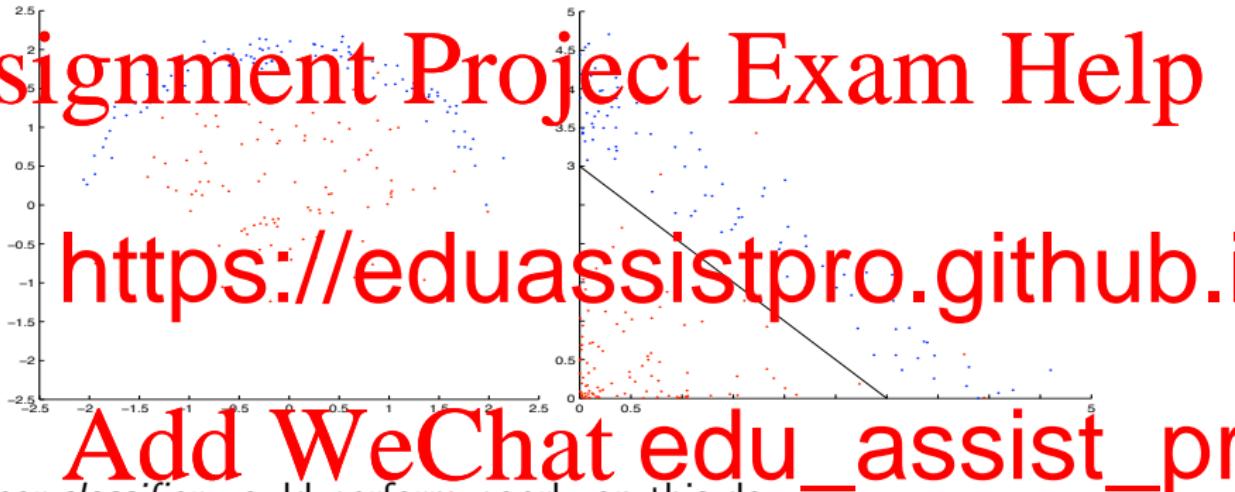
these two f

<https://eduassistpro.github.io>

$$\mathbf{x}'_1 \cdot \mathbf{x}'_2 = x_1^2x_2^2 + y_1^2y_2^2 + 2x_1y_1x_2y_2 = (x_1x_2 + y_1y_2)^2 = (\mathbf{x}_1 \cdot \mathbf{x}_2)^2$$

That is, by squaring the dot product in the original space we get the dot product in the new space *without actually calculating the vectors!* A function that calculates the dot product in feature space directly from the vectors in the original space is called a *kernel* – here the kernel is $\kappa(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1 \cdot \mathbf{x}_2)^2$.

Non-linearly separable data



A *linear classifier* would perform poorly on this data.

Converting the original (x, y) data into $(x', y') = (x^2, y^2)$, the data becomes more ‘linear’, and a linear decision boundary $x' + y' = 3$ separates the data fairly well. In the original space this corresponds to a circle with radius $\sqrt{3}$ around the origin.

Where do features come from ?

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Feature engineering

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Can features be learned ?

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Yes, to some extent. For example, in the intermediate layers of a convolutional neural network¹.

¹Image downloaded 28/2/18 from

<http://devblogs.nvidia.com/deep-learning-nutshell-core-concepts/>

Where does data come from ?

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

AlphaZero played around 44 million chess games against an expert (itself).

Tasks for machine learning

Assignment Project Exam Help

The most common machine learning tasks are *predictive*, in the sense that they conc

- Bin
- Reg
- Clustering: hidden target
- Dimensionality reduction: intrinsic structure

Add WeChat edu_assist_pro

Exploratory or descriptive tasks are concerned with the underlying structure in the data.

Measuring similarity I

Assignment Project Exam Help

If our e-mails are described by word-occurrence features as in the text classification example, the similarity of e-mails would be measured in terms of the word

number of common words occurring in both e-mails (cosine coefficient).

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

Suppose that one e-mail contains 42 (different) words, another e-mail contains 112 words, and the two e-mails have 23 words in common. Then their similarity would be $\frac{23}{42+112-23} = \frac{23}{130} = 0.18$.

Measuring similarity II

Assignment Project Exam Help

We can then cluster our e-mails into groups, such that the average similarity of an e-mail to the other e-mails in its group is much larger than the average

<https://eduassistpro.github.io>

While it would

separated clusters corresponding to spam and ham – this

– the clusters may reveal some interesting and useful structures

It may be possible to identify a particular kind of spam in this

subgroup uses a vocabulary, or language, not found in other

Add WeChat `edu_assist_pro`

Clustering

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

We want to cluster similarly expressed genes in cancer samples.

How many clusters ? |

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

How many clusters ? II

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Looking for structure I

Consider the following matrix:

Assignment Project Exam Help

| | | | |
|---|---|---|---|
| 1 | 0 | 1 | 0 |
| 0 | 2 | 2 | 2 |

<https://eduassistpro.github.io/ml-structure/>

Imagine these represent ratings by six different people of 0 to 3, on four different films – say *The Shawshank Redemption*, *The Usual Suspects*, *The Godfather*, and *The Big Lebowski*.

left to right). *The Godfather* seems to be the most popular of the four with an average rating of 1.5, and *The Shawshank Redemption* is the least appreciated with an average rating of 0.5.

Can you see any structure in this matrix?

Looking for structure II

Assignment Project Exam Help

$$\begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 2 & 2 & 2 \\ 0 & 0 & 0 & 1 \\ 1 & 2 & 3 & 2 \\ 1 & 0 & 1 & 1 \\ 0 & 2 & 2 & 3 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 2 & 0 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

<https://eduassistpro.github.io>

- The right-most matrix associates films (in columns) *The Shawshank Redemption* and *The Godfather* to two different genres, say drama and crime, *The Big Lebowski* belongs to both, and *The Godfather* is a crime and a new genre (say comedy).
- The tall, 6-by-3 matrix then expresses people's preferences in terms of genres.

Looking for structure III

Assignment Project Exam Help

- Fin
imp
pref <https://eduassistpro.github.io>

Add WeChat edu_assist_pro

The philosophical problem

Assignment Project Exam Help

Deduction: derive specific consequences from general theories

Induction

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

Deduction is well-founded (mathematical logic).

Induction is (philosophically) problematic – induct

often seems to work – an inductive argument !

Add WeChat `edu_assist_pro`

Generalisation - the key objective of machine learning

What we are really interested in is *generalising* from the sample of data in our training set. This can be stated as:

The induced function

Any function

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

well approximates the target function well over other examples.

Add WeChat edu_assist_pro

A corollary of this is that it is necessary to make some assumptions about the type of target function in a task for an algorithm to go beyond the data, i.e., generalise or learn.

Cross-validation I

Assignment Project Exam Help

There are c

We use the d

use a separate *validation* or *development* set.

Add WeChat edu_assist_pro

Cross-validation II

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Cross-validation III

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Cross-validation IV

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

Contingency table I

Assignment Project Exam Help

Two-class prediction case:

| | | https://eduassistpro.github.io |
|-----|---------------------|--------------------------------|
| Yes | True Positive (TP) | |
| No | False Positive (FP) | |

Add WeChat `edu_assist_pro`

Contingency table II

Classification Accuracy on a sample of labelled pairs $(x, c(x))$ given a learned classification model that predicts, for each instance x , a class value $\hat{c}(x)$:

$$\frac{1}{|Test|} \sum_{x \in Test} I[\hat{c}(x) = c(x)]$$

where $Test$ is a test set and $I[]$ is the indicator argument evaluates to true, and 0 otherwise.

Classification Error is $1 - \text{acc.}$

Overfitting

Imagine you are preparing for your *Machine Learning 101* exam. Helpfully, your professor has made previous exam papers and their worked answers available online. You begin by trying to answer the questions from previous papers an

Unfortu <https://eduassistpro.github.io/>
the model a
completely consists of past questions, you are certain t
if the new exam asks different questions about the same
would be ill-prepared and get a much lower mark than wi
traditional preparation.

In this case, one could say that you were *overfitting* the past exam papers and that the knowledge gained didn't *generalise* to future exam questions.

The Bias-Variance Decomposition

Assignment Project Exam Help

Error (particularly MSE) can be shown to have two components:

Bias: error due to function being learned

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

Variance: error due to variation between training distribution and data generated by the target function

Add WeChat edu_assist_pro

Inductive Bias

Assignment Project Exam Help

All in

<https://eduassistpro.github.io>

Box & Drape

Add WeChat edu_assist_pro

Inductive Bias

Confusingly, “inductive bias” is *NOT* the same “bias” as in the “bias-variance” decomposition.

Assignment Project Exam Help

“Inducti

on the mod

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

Essentially it means that the algorithm and model com
using to solve the learning problem is appropriate for th

Add WeChat edu_assist_pr

Success in machine learning requires understandin
algorithms and models, and choosing them appropriately for the task².

²Even true for “deep learning”, but watch Andrew Ng’s talk on this
at <http://www.youtube.com/watch?v=F1ka6a13S9I&t=48s>.

No Free Lunch

Assignment Project Exam Help

*Uni
off-*

<https://eduassistpro.github.io>

Wolpert (1996)

Add WeChat edu_assist_pro

No Free Lunch

Assignment Project Exam Help

Owing to the ‘‘No Free Lunch’’ Theorem, there is no universally best machine learning algorithm.

Some learners make assumptions about their learners.
their assumptions are often wrong.
for the learners.

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

On some other tasks, those assumptions, and hence the learners, perform much worse than others.

Add WeChat `edu_assist_pro`

These assumptions form the inductive bias of the learning algorithm.

Ethics of machine learning

Machine learning algorithms are now widely available and increasingly used in “big data” in commercial, medical, legal and scientific applications.

Assignment Project Exam Help

These applications bring many benefits, but as they become more widespread

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

In some of the

machine learning must be considered

Add WeChat edu_assist_pro.com

For example, can the use of machine learning be biased if the training data does not adequately represent all of the data? Can the model be applied? Here is a recent paper discussing some of these issues and giving some recommendations³.

³Zook et al. (2017) Ten simple rules for responsible big data research. PLoS Comput Biol 13(3): e1005399. <http://doi.org/10.1371/journal.pcbi.1005399>.

Summary

The purpose of this introductory lecture was, from a high-level, to survey the landscape that we will explore in this course.

Assignment Project Exam Help

We have motivated the subject matter, outlined what parts of machine learning we will cover, and termi

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

Through

programming languages, machine learning tools and frameworks.

These change very rapidly, although the core techniques remain. Challenges remain.

Add WeChat edu_assist_pro

Following this lecture you should have a clear idea of the course scope. In the remaining lectures we will expand on the topics we have just introduced and go into more detail. You will also get to work on practical applications of what we have covered.