



Statistical Machine Learning

Assignment Project Exam Help

<https://eduassistpro.github.io>

Canberra
Semester One, 2021.
Add WeChat [edu_assist_pro](#)

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



Assignment Project Exam Help

Part I

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr



Assignment Project Exam Help

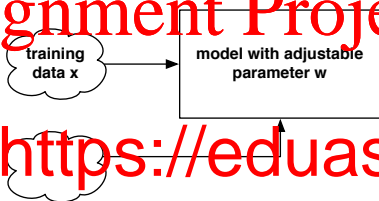
- Bas
- Ma
- Re
- Bias variance decomposition

<https://eduassistpro.github.io>

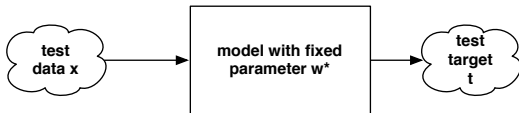
Add WeChat edu_assist_pr



Training Phase



Test Phase



Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



Assignment Project Exam Help

- Bayes Theorem

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{p(\mathbf{t} | \mathbf{w})} \quad p(\mathbf{t} | \mathbf{w}) = \frac{p(\mathbf{t} | \mathbf{w}) p(\mathbf{w})}{p(\mathbf{t})}$$

- wh
and
• l.i.d

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

$$\begin{aligned} p(\mathbf{t} | \mathbf{w}) &= \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^\top \phi(\mathbf{x}_n), \beta^{-1}) \\ &= \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^\top \phi(\mathbf{x}_n), \beta^{-1}) \\ &= \text{const} \times \exp\left\{-\beta \frac{1}{2} (\mathbf{t} - \Phi \mathbf{w})^\top (\mathbf{t} - \Phi \mathbf{w})\right\} \\ &= \mathcal{N}(\mathbf{t} | \Phi \mathbf{w}, \beta^{-1} \mathbf{I}) \end{aligned}$$

How to choose a prior?



- The choice of prior affords an intuitive control over our inductive bias
- All inference schemes have such biases, and often arise more
- Can

<https://eduassistpro.github.io>

An answer to the second question:

Definition (Conjugate Prior)

A class of prior probability distributions $p(w)$ such that for any class of likelihood functions $p(x | w)$ if the resulting posterior distributions $p(w | x)$ are in the same family as $p(w)$.

Examples of Conjugate Prior Distributions



Table: Discrete likelihood distributions

Likelihood	Conjugate Prior
Bernoulli	Beta

<https://eduassistpro.github.io>

Table: Continuous Likelihood distributions

Likelihood	Conjugate
Uniform	Pareto
Exponential	Gamma
Normal	Normal (mean parameter)
Multivariate normal	Multivariate normal (mean parameter)

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

Conjugate Prior to a Gaussian Distribution



- Example : If the likelihood function is Gaussian, choosing a Gaussian prior for the mean will ensure that the posterior distribution is also Gaussian.
- Given a marginal distribution for \mathbf{x} and a conditional Gaussian distribution for \mathbf{y} given \mathbf{x} in the form

<https://eduassistpro.github.io>

- we get

Add WeChat edu_assist_pro

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1}) +$$
$$p(\mathbf{x} | \mathbf{y}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\Sigma}\{\mathbf{A}^\top \mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma})$$

where $\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^\top \mathbf{L} \mathbf{A})^{-1}$.

Note that the covariance $\boldsymbol{\Sigma}$ does not involve \mathbf{y} .



Conjugate Prior to a Gaussian Distribution (intuition)

Given

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$

$$p(\mathbf{y} | \mathbf{x}) = \mathcal{N}(\mathbf{y} | \mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}) \Leftrightarrow \mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b} + \mathcal{N}(\mathbf{0}, \mathbf{L}^{-1})$$

We have $\mathbb{E}[\mathbf{y}] = \mathbb{E}[\mathbf{A}\mathbf{x} + \mathbf{b}] = \mathbf{A}\boldsymbol{\mu} + \mathbf{b}$ and by the easily proven Bienaymé formula for t

$$\text{cov}[\mathbf{y}] = \text{cov}[\mathbf{A}\mathbf{x} + \mathbf{b}] + \text{cov}[\mathbf{L}^{-1}]$$

So \mathbf{y} is Gau

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1})$$

Then letting $\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^\top \mathbf{L} \mathbf{A})^{-1}$ and

$$p(\mathbf{x} | \mathbf{y}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\Sigma} \{\mathbf{A}^\top \mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda} \boldsymbol{\mu}\})$$

$$\Leftrightarrow \mathbf{x} = \boldsymbol{\Sigma} \{\mathbf{A}^\top \mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda} \boldsymbol{\mu}\} + \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$$

yields the correct moments for \mathbf{x} , since

$$\begin{aligned} \mathbb{E}[\mathbf{x}] &= \mathbb{E}[\boldsymbol{\Sigma} \{\mathbf{A}^\top \mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda} \boldsymbol{\mu}\}] = \boldsymbol{\Sigma} \{\mathbf{A}^\top \mathbf{L}(\mathbf{A}\boldsymbol{\mu} + \mathbf{b} - \mathbf{b}) + \boldsymbol{\Lambda} \boldsymbol{\mu}\} \\ &= \boldsymbol{\Sigma} \{\mathbf{A}^\top \mathbf{L} \mathbf{A} \boldsymbol{\mu} + \boldsymbol{\Lambda} \boldsymbol{\mu}\} = (\boldsymbol{\Lambda} + \mathbf{A}^\top \mathbf{L} \mathbf{A})^{-1} \{\mathbf{A}^\top \mathbf{L} \mathbf{A} + \boldsymbol{\Lambda}\} \boldsymbol{\mu} = \boldsymbol{\mu}, \end{aligned}$$

and it is similar (but tedious ; don't do it) to recover $\text{cov}[\mathbf{x}] = \boldsymbol{\Lambda}$.

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



- Choose a Gaussian prior with mean \mathbf{m}_0 and covariance \mathbf{S}_0

Assignment Project Exam Help

- Same likelihood as before (here written in vector form):

- Giv

<https://eduassistpro.github.io>

$$p(\mathbf{w} | \mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)$$

where

Add WeChat edu_assist_pr

$$\mathbf{m}_N = \mathbf{S}_N(\mathbf{S}_0^{-1}\mathbf{m}_0 + \beta\Phi^T\mathbf{t})$$

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta\Phi^T\Phi$$

(derive this with the identities on the previous slides)



- For simplicity we proceed with $\mathbf{m}_0 = 0$ and $\mathbf{S}_0 = \alpha^{-1}\mathbf{I}$, so

Assignment Project Exam Help

- The posterior becomes $p(\mathbf{w} | \mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)$ with

<https://eduassistpro.github.io>

- For $\alpha \ll \beta$ we get

Add WeChat [edu_assist_pro](#)

- Log of posterior is sum of log likelihood and log of prior

$$\ln p(\mathbf{w} | \mathbf{t}) = -\frac{\beta}{2}(\mathbf{t} - \Phi\mathbf{w})^\top (\mathbf{t} - \Phi\mathbf{w}) - \frac{\alpha}{2}\mathbf{w}^\top \mathbf{w} + \text{const}$$



- Log of posterior is sum of log likelihood and log of prior

$$\log p(\mathbf{w} | \mathbf{t}) = -\beta \underbrace{\frac{1}{2} \|\mathbf{t} - \mathbf{S}\mathbf{w}\|^2}_{\text{sum-of-squares-error}} - \underbrace{\frac{\alpha}{2} \|\mathbf{w}\|^2}_{\text{regulariser}} + \text{const.}$$

- The

<https://eduassistpro.github.io>

corresponds to minimising the sum-of-squares function with quadratic regularisation coefficient

- The posterior is Gaussian so mode =
- For $\alpha \ll \beta$ we recover unregularised least squares (equivalently m.a.p. approaches maximum likelihood), for example in case of
 - an infinitely broad prior with $\alpha \rightarrow 0$
 - an infinitely precise likelihood with $\beta \rightarrow \infty$

Bayesian Inference in General: Sequential Update of Belief

Statistical Machine
Learning

© 2021

Ong & Walder & Webers
& Xie

Data61 | CSIRO
ANU Computer Science



Assignment Project Exam Help

- If we have not yet seen any datapoint ($N = 0$), the posterior is equal to the prior.

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

- Sequential arrival of data points : the posterior given some observed data acts as the prior for the future data.
- Nicely fits a sequential learning framework.



Assignment Project Exam Help

- Example of a linear basis function model

- Single input x , single output t

- Lin

- Tru

1

2 Calculate $f(x_n, \mathbf{a}) = a_0 + a_1 x_n$, where $a_0 = 0.3$, $a_1 = 0.5$.

3 Add Gaussian noise with standard deviation

Add WeChat: edu_assist_pr

- Set the precision of the uniform prior to $\alpha = 2.0$.



Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr



Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr



Assignment Project Exam Help

- In the training phase, data \mathbf{x} and targets \mathbf{t} are provided
- In the test phase, a new data value x is given and the
- Ba
giv
targets \mathbf{t}

<https://eduassistpro.github.io>

$$p(t|x, \mathbf{x}, \mathbf{t})$$

- This is the Predictive Distribution (c.f. t distribution, which is over the parameters).

Add WeChat edu_assist_pr

How to calculate the Predictive Distribution?



- Introduce the model parameter \mathbf{w} via the sum rule

$$\begin{aligned} p(t|x, \mathbf{x}, \mathbf{t}) &= \int p(t, \mathbf{w} | x, \mathbf{x}, \mathbf{t}) d\mathbf{w} \\ &= \int p(t | \mathbf{w}, x, \mathbf{x}, \mathbf{t}) p(\mathbf{w} | x, \mathbf{x}, \mathbf{t}) d\mathbf{w} \end{aligned}$$

- The t
mo
trai

<https://eduassistpro.github.io>

$$p(t | \mathbf{w}, x, \mathbf{x}, \mathbf{t}) = p(t | \mathbf{w})$$

- The model parameter \mathbf{w} are learned with \mathbf{x} and the training targets \mathbf{t} only

$$p(\mathbf{w} | x, \mathbf{x}, \mathbf{t}) = p(\mathbf{w} | \mathbf{x}, \mathbf{t})$$

- Predictive Distribution

$$p(t|x, \mathbf{x}, \mathbf{t}) = \int p(t | \mathbf{w}, x) p(\mathbf{w} | \mathbf{x}, \mathbf{t}) d\mathbf{w}$$

Add WeChat [edu_assist_pro](#)

Proof of the Predictive Distribution



The predictive distribution is

$$p(t|x, \mathbf{x}, \mathbf{t}) = \int p(t|\mathbf{w}, x, \mathbf{x}, \mathbf{t})p(\mathbf{w}|x, \mathbf{x}, \mathbf{t})d\mathbf{w}$$

because

$$\int p(t|$$

<https://eduassistpro.github.io>

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

or simply

$$\begin{aligned} \int p(t|\mathbf{w}, x, \mathbf{x}, \mathbf{t})p(\mathbf{w}|x, \mathbf{x}, \mathbf{t})d\mathbf{w} &= \int p(t, \mathbf{w}|x, \mathbf{x}, \mathbf{t})d\mathbf{w} \\ &= p(t|x, \mathbf{x}, \mathbf{t}). \end{aligned}$$



- Find the predictive distribution

$$p(t|\mathbf{t}, \alpha, \beta) = \int p(t|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w}$$

(remember : conditioning on \mathbf{x} is often suppressed to
sim

- No
on

$$p(t|\mathbf{w}, \beta) = \mathcal{N}(t|\mathbf{w}^\top \phi(\mathbf{x}), \beta^{-1})$$

- and the posterior was

$$p(\mathbf{w}|\mathbf{t}, \alpha, \beta) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N)$$

where

$$\mathbf{m}_N = \beta \mathbf{S}_N \Phi^\top \mathbf{t}$$

$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^\top \Phi$$



- If we do the integral (it turns out to be the convolution of the two Gaussians), we get for the predictive distribution

$$p(t | \mathbf{x}, \mathbf{t}, \alpha, \beta) \equiv \mathcal{N}(t | \mathbf{m}_N^\top \phi(\mathbf{x}), \sigma_N^2(\mathbf{x}))$$

wh

<https://eduassistpro.github.io>

- This is more easily shown using a similar approach earlier “intuition” slide and again with the Binary formula, now using

$$t = \mathbf{w}^\top \phi(\mathbf{x}) + \mathcal{N}(0, \beta^{-1}).$$

However this is a linear-Gaussian specific trick and in general we need to integrate out the parameters.

Predictive Distribution with Simplified Prior

Statistical Machine
Learning

© 2021

Ong & Walder & Webers
& Xie

Data61 | CSIRO
ANU Computer Science



Example with artificial sinusoidal data from $\sin(2\pi x)$ (green)
and added noise. Number of data points $N = 1$.



Mean of the predictive distribution (red) and regions of one
standard deviation from mean (red shaded).

Predictive Distribution with Simplified Prior

Statistical Machine
Learning

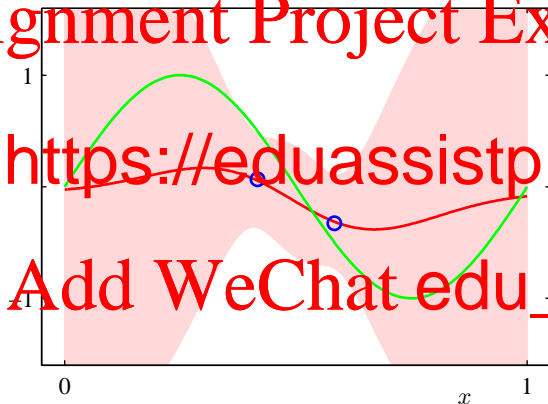
© 2021

Ong & Walder & Webers
& Xie

Data61 | CSIRO
ANU Computer Science



Example with artificial sinusoidal data from $\sin(2\pi x)$ (green)
and added noise. Number of data points $N = 2$.



Mean of the predictive distribution (red) and regions of one
standard deviation from mean (red shaded).

Predictive Distribution with Simplified Prior

Statistical Machine
Learning

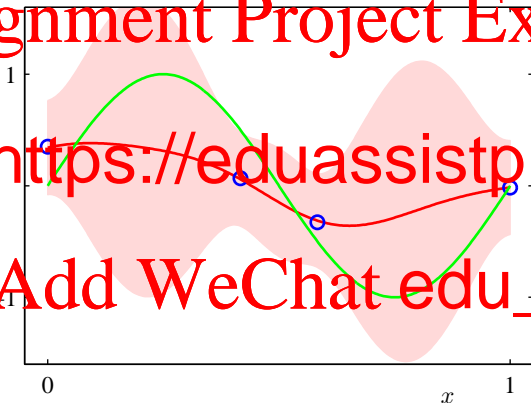
© 2021

Ong & Walder & Webers
& Xie

Data61 | CSIRO
ANU Computer Science



Example with artificial sinusoidal data from $\sin(2\pi x)$ (green)
and added noise. Number of data points $N = 4$.



Mean of the predictive distribution (red) and regions of one
standard deviation from mean (red shaded).

Predictive Distribution with Simplified Prior

Statistical Machine
Learning

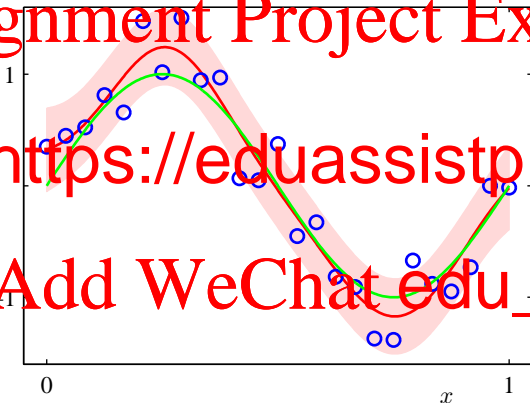
© 2021

Ong & Walder & Webers
& Xie

Data61 | CSIRO
ANU Computer Science



Example with artificial sinusoidal data from $\sin(2\pi x)$ (green)
and added noise. Number of data points $N = 25$.

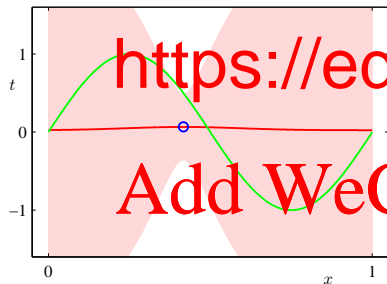


Mean of the predictive distribution (red) and regions of one
standard deviation from mean (red shaded).



Assignment Project Exam Help

Plots of the function $y(x; \mathbf{w})$ using samples from the posterior distribution over \mathbf{w} . Number of data points $N = 1$.



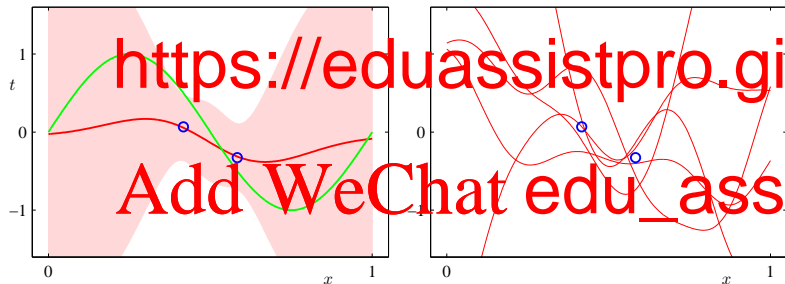
<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



Assignment Project Exam Help

Plots of the function $y(x; \mathbf{w})$ using samples from the posterior distribution over \mathbf{w} . Number of data points $N = 2$.



<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Predictive Distribution with Simplified Prior

Statistical Machine
Learning

© 2021

Ong & Walder & Webers
& Xie

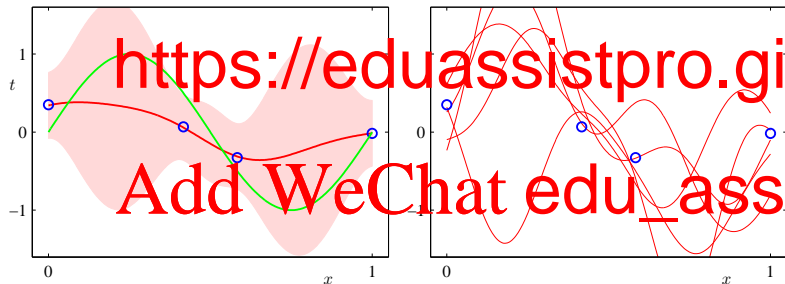
Data61 | CSIRO

ANU Computer Science



Assignment Project Exam Help

Plots of the function $y(x; \mathbf{w})$ using samples from the posterior distribution over \mathbf{w} . Number of data points $N = 4$.



<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

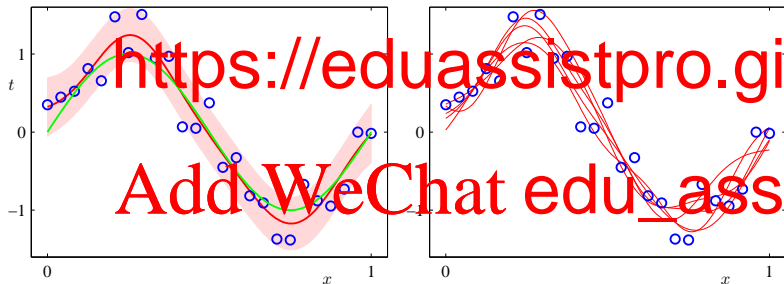
Predictive Distribution with Simplified Prior

Statistical Machine
Learning

© 2021
Ong & Walder & Webers
& Xie
Data61 | CSIRO
ANU Computer Science



Plots of the function $y(x; w)$ using samples from the posterior distribution over w . Number of data points $N = 25$.



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

Limitations of Linear Basis Function Models

Statistical Machine
Learning

© 2021

Ong & Walder & Webers
& Xie

Data61 | CSIRO

ANU Computer Science



- Basis function $\phi_j(\mathbf{x})$ are fixed before the training data set is observed

- Curse of dimensionality : Number of basis function grows rapidly

- But this can be avoided

- dimension much smaller than D . Need algorithms which place basis functions only where data are (e.g. methods / Gaussian processes)
- Target variables may only depend on a few significant directions within the data manifold. Need algorithms which can exploit this property (e.g. linear methods or shallow neural networks).

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



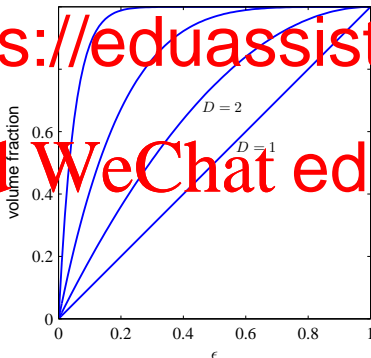
- Linear Algebra allows us to operate in n -dimensional vector spaces using the intuition from our 3-dimensional world as a vector space. No surprises as long as n is finite.
- If we add more structure to a vector space (e.g. inner product, 3-dimensional volume), the volume of the sphere in a D -dimensional space which lies between radius $r = 1$ and $r = 1 - \epsilon$
- Volume scales like r^D , therefore the formula of a sphere is $V_D(r) = K_D r^D$.

$$\frac{V_D(1) - V_D(1 - \epsilon)}{V_D(1)} = 1 - (1 - \epsilon)^D$$



- Fraction of the volume of the sphere in a D -dimensional space which lies between radius $r = 1$ and $r = 1 - \epsilon$

Assignment Project Exam Help

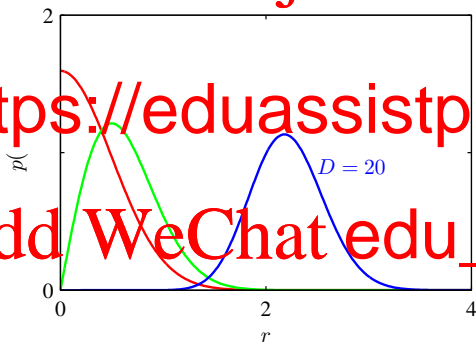
$$\frac{V_D(1) - V_D(1 - \epsilon)}{V_D(1)} = 1 - (1 - \epsilon)^D$$


<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



- Probability density with respect to radius r of a Gaussian distribution for various values of the dimensionality D



Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



- Probability density with respect to radius r of a Gaussian distribution for various values of the dimensionality D .

- Example: $D = 2$; assume $\mu = 0, \Sigma = I$

$$\mathcal{N}(x | 0, I) = \frac{1}{2\pi} \exp\left(-\frac{1}{2}x^\top x\right) = \frac{1}{2\pi} \exp\left(-\frac{1}{2}(x_1^2 + x_2^2)\right)$$

- Co

<https://eduassistpro.github.io>

- Probability in the new coordinates

$$p(r, \phi | 0, I) = \mathcal{N}(r(x), \phi(x) | 0, I)$$

where $|J| = r$ is the determinant of the Jacobian for the given coordinate transformation.

$$p(r, \phi | 0, I) = \frac{1}{2\pi} r \exp\left\{-\frac{1}{2}r^2\right\}$$



- Probability density with respect to radius r of a Gaussian distribution for $D = 2$ (and $\mu = 0, \Sigma = I$)

Assignment Project Exam Help

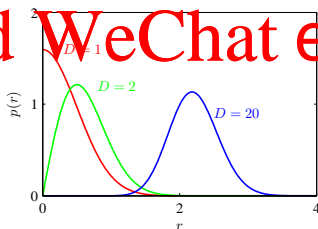
$$p(r; \mu=0, I) = \frac{1}{2\pi} r \exp\left\{-\frac{r^2}{2}\right\}$$

- Inte

<https://eduassistpro.github.io>

0

Add WeChat edu_assist_pr





Assignment Project Exam Help

- Basis functions
- Maximum likelihood with Gaussian noise
- Re
- Bia
- Co
- Bayesian linear regression
- Sequential update of the posterior
- Predictive distribution
- Curse of dimensionality

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr