



Outlines

Overview

Introduction

Linear Algebra

Probability

Linear Regression

Linear Regression

Linear Classification 1

Linear Classification 2

Kernel Methods

Sparse Kernel Methods

Structure Models and EM 1

Structure Models and EM 2

Neural Networks 1

Neural Networks 2

Deep Component Analysis

Generators

Graphical Models 1

Graphical Models 2

Graphical Models 3

Sampling

Sequential Data 1

Sequential Data 2

Statistical Machine Learning

Assignment Project Exam Help

Christian Walder

<https://eduassistpro.github.io>

College of Engineering and Computer Science
The Australian National University

Add WeChat edu_assist_pro

Canberra

Semester One, 2020.

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



Assignment Project Exam Help

Part III

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Sparse

Bias-Variance
Composition



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Sparse

Bias-Variance
Composition

Assignment Project Exam Help

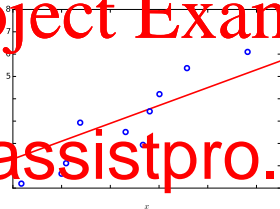
$N = 10$

$\mathbf{x} \equiv (x_1, \dots, x_N)^T$

\mathbf{t}

x_i

t_i



<https://eduassistpro.github.io>

- Predictor $y(x, \mathbf{w})$?
- Performance measure?
- Optimal solution \mathbf{w}^* ?
- Recall: projection, inverse

Add WeChat edu_assist_pro



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Sparse

Bias-Variance
Composition

Assignment Project Exam Help

- Ga
- Ba
- Exp

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



Review

Linear Basis Function Models

Maximum Likelihood and Least Squares

Squared Loss

Regularized Least Squares

Multiple Outputs

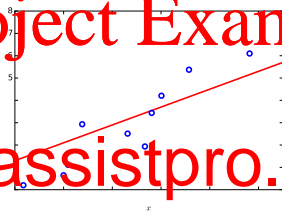
Selection for

Bias-Variance
Composition

$$N = 10$$

$$\mathbf{x} \equiv (x_1, \dots, x_N)^T$$

$$\mathbf{t} \equiv (t_1, \dots, t_N)^T$$



$y(x,$

$$X \equiv [\mathbf{x} \quad \mathbf{1}]$$

$$\mathbf{w}^* = (X^T X)^{-1} X^T \mathbf{t}$$

Add WeChat edu_assist_pro

We assume

$$t = \underbrace{y(\mathbf{x}, \mathbf{w})}_{\text{deterministic}} + \underbrace{\epsilon}_{\text{Gaussian noise}}$$



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Prediction

Bias-Variance
Composition

Assignment Project Exam Help

- a *prior* belief about the parameter \mathbf{w} captured in the prior probability $p(\mathbf{w})$

- obs

- cal
obs

<https://eduassistpro.github.io>

$$\frac{1}{p(\mathcal{D})}$$

- $p(\mathcal{D} | \mathbf{w})$ as a function of \mathbf{w} . **likelihood fun**
- likelihood expresses how probable the data are for different values of \mathbf{w} — it is **not** a probability d respect \mathbf{w} (but it is with respect to \mathcal{D} ; prove it)

Add WeChat: edu_assist_pro



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Regression

Bias-Variance
Composition

Assignment Project Exam Help

- Consider the linear regression problem, where we have random variables \mathbf{x}_n and t_n .

- We

- We

<https://eduassistpro.github.io>

For a given θ the density defines the probability of observing t_n given \mathbf{x}_n .

- We are interested in finding θ that **maximize** the probability (called the **likelihood**) of the data.

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



Review

Linear Basis Function Models

Maximum Likelihood and Least Squares

Squared Loss

Regularized Least Squares

Multiple Outputs

Selection for

Bias-Variance Composition

Assignment Project Exam Help

<https://eduassistpro.github.io>

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

Frequentist Approach

- Likelihood function $p(\mathcal{D}|\theta, w)$
- v some 'estimator'
- error bars on the estimated v obtained from the distribution of possible data sets \mathcal{D}

Bayesian Approach

- prior
- from data distribution



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Prediction

Bias-Variance
Composition

Assignment Project Exam Help

- choose w for which the likelihood $p(\mathcal{D} | w)$ (the probability of the observed data) is maximal
- the most common heuristic for learning a single fixed w
- equ
fun
- log i
- maximising the likelihood \iff minimi
- Example: Fair-looking coin is tossed three times
landing on heads.
- Maximum likelihood estimate of the probability of landing
heads will give 1.

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Prediction

Bias-Variance
Composition

- including prior knowledge easy (via prior w)
- subjective choice of prior, allows better results by incorporating domain knowledge
- sometimes choice of prior motivated by convenient mathematical form
- prior
- need
 - advances in approximation schemes (Variational Expectation Propagation)
- there is no true w .

Assignment Project Exam Help
<https://eduassistpro.github.io>
Add WeChat edu_assist_pro



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

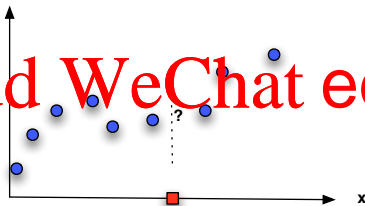
Selection for
Regression

Bias-Variance
Composition

- Given a training data set of N observations $\{\mathbf{x}_n\}$ and target values t_n .

- Goal: Learn to predict the value of one or more target values t given a new value of the input \mathbf{x} .

- Example



Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Supervised Learning: (non-Bayesian) Point Estimate



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

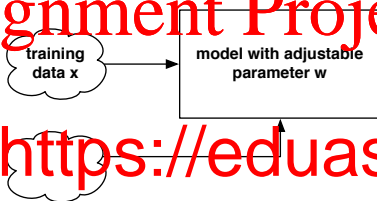
Regularized Least
Squares

Multiple Outputs

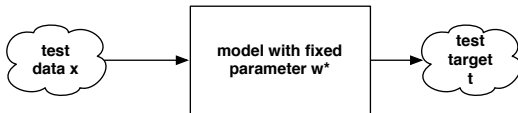
Selection for
Prediction

Bias-Variance
Composition

Training Phase



Test Phase



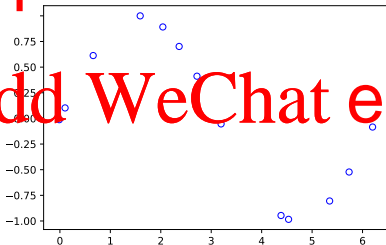
Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Why Linear Regression?

- Analytic solution when minimising sum of squared errors
- Well understood statistical behaviour
- Efficient algorithms exist for convex losses and regularizers
- But





Assignment Project Exam Help

- Linear combination of fixed nonlinear basis functions

<https://eduassistpro.github.io>

- parameter $\mathbf{w} = (w_0, \dots, w_{M-1})$
- basis functions $\phi(\mathbf{x}) = (\phi_0(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x}))$
- convention $\phi_0(\mathbf{x}) = 1$
- w_0 is the bias parameter

Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Sparse

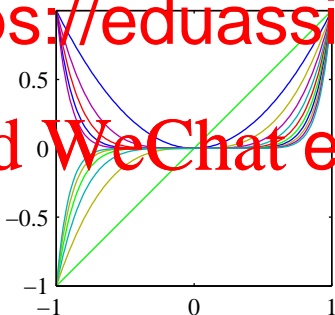
Bias-Variance
Composition



- Scalar input variable x

$$\phi_j(x) = x^j$$

- Limitation: Polynomials are global functions of the input variable x so the learned function will extrapolate poorly



<https://eduassistpro.github.io>

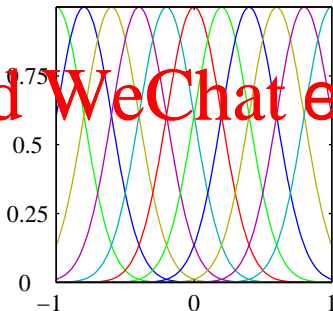
Add WeChat edu_assist_pro

'Gaussian' Basis Functions

- Scalar input variable x

$$\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\}$$

- Not a probability distribution.
- No normalisation required, taken care of by the model par
- We



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Model

Bias-Variance
Composition

Sigmoidal Basis Functions

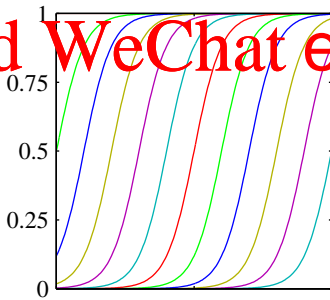
- Scalar input variable x

$$\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s_j}\right)$$

where $\sigma(a)$ is the logistic sigmoid function defined by

$$\frac{1}{1 + e^{-a}}$$

- $\sigma(a)$
tan

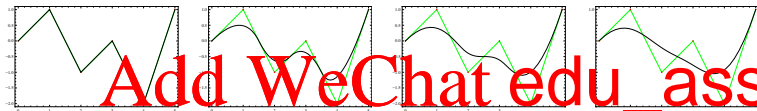


Other Basis Functions

- Fourier Basis : each basis function represents a specific frequency and has infinite spatial extent.
- Wavelets : localised in both space and frequency (also mutually orthogonal to simplify application).

• Splines

<https://eduassistpro.github.io>



Linear
Splines

Quadratic
Splines

Cubic
Splines

Splines

Approximate the points
 $\{(0, 0), (1, 1), (2, -1), (3, 0), (4, -2), (5, 1)\}$ by different
splines.

Maximum Likelihood and Least Squares

Statistical Machine
Learning

© 2020

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



- No special assumption about the basis functions $\phi_j(\mathbf{x})$. In the simplest case, one can think of $\phi_j(\mathbf{x}) = x_j$, or $\phi(\mathbf{x}) = \mathbf{x}$.

- Assume target t is given by

$$t = \underline{y}(\underline{\mathbf{x}}, \underline{\mathbf{w}}) + \epsilon$$

Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

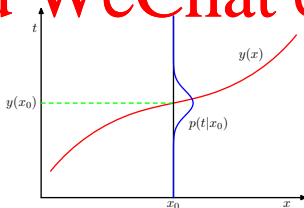
Selection for
Sparse Models

Bias-Variance
Composition

wh
pre

- Thus

$$p(t | \underline{\mathbf{x}}, \underline{\mathbf{w}}, \beta) = \mathcal{N}(t | y(\underline{\mathbf{x}}, \underline{\mathbf{w}}), \beta)$$



Add WeChat edu_assist_pr

Maximum Likelihood and Least Squares



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Sparse Models

Bias-Variance
Composition

- Likelihood of one target t given the data \mathbf{x} was

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

- Now, a set of inputs \mathbf{X} with corresponding target values \mathbf{t} .

- Ass

(i.i.

sa

<https://eduassistpro.github.io>

$$\begin{aligned} p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) &= \prod_{n=1}^N \mathcal{N}(t_n | y(\mathbf{x}_n, \mathbf{w}), \beta^{-1}) \\ &= \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \end{aligned}$$

- From now on drop the conditioning variable \mathbf{X} from the notation, as with supervised learning we do not seek to model the distribution of the input data.

Maximum Likelihood and Least Squares

Statistical Machine
Learning

© 2020

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



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Square and Linear

Regularized Least
Squares

Multiple Outputs

Selection for
Regression

Bias-Variance
Composition

- Consider the **logarithm of the likelihood** $p(\mathbf{t} | \mathbf{w}, \beta)$ (the logarithm is a monotone function!)

$$\ln p(\mathbf{t} | \mathbf{w}, \beta) = \sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1})$$

<https://eduassistpro.github.io>

$$= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \frac{\beta}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2$$

where the **sum-of-squares error function** is

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2.$$

- $\arg \max_{\mathbf{w}} \ln p(\mathbf{t} | \mathbf{w}, \beta) \rightarrow \arg \min_{\mathbf{w}} E_D(\mathbf{w})$

Add WeChat [edu_assist_pro](#)



Review

Linear Basis Function Models

Maximum Likelihood and Least Squares

Squared Loss

Regularized Least Squares

Multiple Outputs

Selection for sparsity

Bias-Variance decomposition

- Goal: Find a more compact representation.

- Rewrite the error function

N

wh

<https://eduassistpro.github.io>

Add WeChat [edu_assist_pro](#)

$$\Phi = \begin{bmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_M(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \dots & \phi_M(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \dots & \phi_M(\mathbf{x}_N) \end{bmatrix}$$

Maximum Likelihood and Least Squares



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Sparse

Bias-Variance
Composition

- The log likelihood is now

$$\begin{aligned}\ln p(\mathbf{t} | \mathbf{w}, \beta) &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w}) \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta \frac{1}{2} (\mathbf{t} - \Phi \mathbf{w})^T (\mathbf{t} - \Phi \mathbf{w})\end{aligned}$$

- Find
- The

$$\nabla_{\mathbf{w}} \ln p(\mathbf{t} | \mathbf{w}, \beta) = \beta \Phi^T (\mathbf{t} - \Phi \mathbf{w}).$$

Setting the gradient to zero gives

$$0 = \Phi^T \mathbf{t} - \Phi^T \Phi \mathbf{w},$$

- which results in

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} = \Phi^\dagger \mathbf{t}$$

where Φ^\dagger is the Moore-Penrose pseudo-inverse of the matrix Φ .



Review

Linear Basis Function Models

Maximum Likelihood and Least Squares

Squashed Linear Regression

Regularized Least Squares

Multiple Outputs

Selection for Regression

Bias-Variance Composition

Maximum Likelihood and Least Squares

- The log likelihood with the optimal \mathbf{w}_{ML} is now

$$\ln p(\mathbf{t} | \mathbf{w}_{ML}, \beta)$$

$$= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \frac{1}{2} (\mathbf{t} - \Phi \mathbf{w}_{ML})^T (\mathbf{t} - \Phi \mathbf{w}_{ML})$$

- Fin

<https://eduassistpro.github.io>

results in

Add WeChat edu_assist_pro

$$\frac{1}{\beta_{ML}} = \frac{1}{N} (\mathbf{t} - \Phi \mathbf{w}_{ML})^T (\mathbf{t} - \Phi \mathbf{w}_{ML})$$

- Note: We can first find the maximum likelihood for \mathbf{w} as this does **not depend** on β . Then we can use \mathbf{w}_{ML} to find the maximum likelihood solution for β .
- Could we have chosen optimisation wrt β first, and then wrt to \mathbf{w} ?

Sequential Learning - Stochastic Gradient Descent

Statistical Machine Learning

© 2020

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



Review

Linear Basis Function Models

Maximum Likelihood and Least Squares

Sequential Learning

Regularized Least Squares

Multiple Outputs

Selection for

Bias-Variance Composition

Assignment Project Exam Help

- For large data sets, calculating the maximum likelihood parameters \mathbf{w}_{ML} and β_{ML} may be costly.

- For
- Us
- If the then

- 1 initialise $\mathbf{w}^{(0)}$ to some starting value
- 2 update the parameter vector at iteration

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta$$

where E_n is the error function after presenting the n th data set, and η is the **learning rate**.

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



Review

Linear Basis Function Models

Maximum Likelihood and Least Squares

Sequential Learning

Regularized Least Squares

Multiple Outputs

Selection for

Bias-Variance Composition

Assignment Project Exam Help

- For the test set

<https://eduassistpro.github.io>

- The value for the learning rate must be chosen carefully. A **too large** learning rate may prevent the algorithm converging. A **too small** learning rate does converge too slowly.

Add WeChat [edu_assist_pro](#)



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Regression

Bias-Variance
Composition

Assignment Project Exam Help

- Add regularisation in order to prevent overfitting

$$E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$$

- Sim

<https://eduassistpro.github.io>

$$E_W(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

- Maximum likelihood solution

Add WeChat edu_assist_pro

$$\mathbf{w} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared L-norming

Regularized Least
Squares

Multiple Outputs

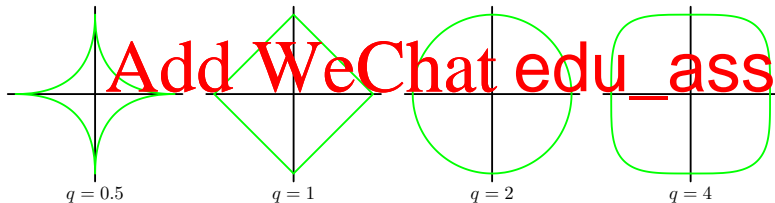
Selection for
Regression

Bias-Variance
Composition

- More general regulariser

$$E_W(\mathbf{w}) = \frac{1}{2} \sum_j |w_j|^q$$

- $q =$



Add WeChat edu_assist_pro



Assignment Project Exam Help

Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared L2 Regulariser

Regularized Least
Squares

Multiple Outputs

Selection for
Sparse

Bias-Variance
Composition

- By the Lagrange multiplier method, minimization of the regularized error function

$$\frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 + \frac{\lambda}{2} \sum_{j=1}^M w_j^2$$

<https://eduassistpro.github.io>

- is eq
sum-of-squares error,

Add WeChat edu_assist_pro

- This yields the figures on the next slide.

Comparison of Quadratic and Lasso Regulariser

Statistical Machine
Learning

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



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Quadratic Learning

Regularized Least
Squares

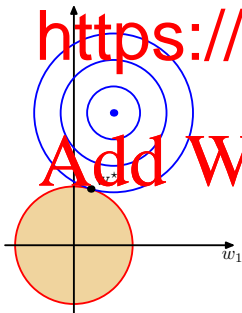
Multiple Outputs

Selection for
Regression

Bias-Variance
Composition

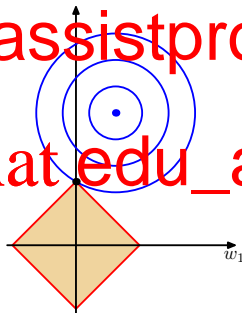
Quadratic regulariser

$$\frac{1}{2} \sum_{j=1}^M w_j^2$$



Lasso regulariser

$$\frac{1}{2} \sum_{j=1}^M |w_j|$$



Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Regression

Bias-Variance
Composition

- More than 1 target variable per data point.
 - \mathbf{y} becomes a vector instead of a scalar. Each dimension can be treated with a different set of basis functions (and that dim

- Here

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

where \mathbf{y} is a K -dimensional column vector,
matrix of model parameters, and

$\phi(\mathbf{x}) = (\phi_0(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x}))$, with $\phi_0(\mathbf{x}) = 1$, as before.

- Define target matrix \mathbf{T} containing the target vector \mathbf{t}_n^T in the n^{th} row.



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Regression

Bias-Variance
Composition

Assignment Project Exam Help

- Suppose the conditional distribution of the target vector is an isotropic Gaussian of the form

- Then

$$\begin{aligned}\ln p(\mathbf{T} | \mathbf{X}, \mathbf{W}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(\mathbf{t}_n | \mathbf{W}^T \phi(\mathbf{x}_n)) \\ &= \frac{NK}{2} \ln \left(\frac{\beta}{2\pi} \right) - \frac{\beta}{2} \sum_{n=1}^N \mathbf{t}_n^T \phi(\mathbf{x}_n) \mathbf{t}_n\end{aligned}$$

Add WeChat edu_assist_pro



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Regression

Bias-Variance
Composition

- Maximisation with respect to \mathbf{W} results in

$$\mathbf{W}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{r}$$

- For

- The
decouples.
- Holds also for a general Gaussian noise distrib
arbitrary covariance matrix.
- Why? \mathbf{W} defines the mean of the Gaussian noise
distribution. And the maximum likelihood solution for the
mean of a multivariate Gaussian is independent of the
covariance.

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Loss Function for
Regression

Bias-Variance
Composition

Assignment Project Exam Help

- Overfitting and underfitting
- Regression correlation
- Frequentist's viewpoint of the model complexity bias-variance trade-off.

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Loss Function for Regression



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Loss Function for
Regression

Bias-Variance
Composition

- Choose an estimator $y(\mathbf{x})$ to estimate the target value t for each input \mathbf{x} .
- Choose a loss function $L(t, y(\mathbf{x}))$ which measures the difference between the target t and the estimate $y(\mathbf{x})$.
- The

<https://eduassistpro.github.io>

- Common choice: Squared Loss

Add WeChat edu_assist_pro

$$L(t, y(\mathbf{x})) = \{y(\mathbf{x}) - t\}^2$$

- Expected loss for squared loss function

$$\mathbb{E}[L] = \iint \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) \, d\mathbf{x} \, dt.$$



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Loss Function for
Regression

Bias-Variance
Composition

Assignment Project Exam Help

2

- Min

$$y(\mathbf{x}) = \frac{\int t p(\mathbf{x}, t) dt}{p(\mathbf{x})} = \int t p(t | \mathbf{x})$$

(calculus of variations is not required to derive this, we may work point-wise by fixing an \mathbf{x} stationarity to solve for $y(\mathbf{x})$ — why is that sufficient?).

Optimal Predictor for Squared Loss

Statistical Machine
Learning

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



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

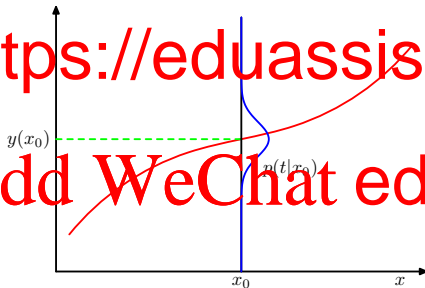
Regularized Least
Squares

Multiple Outputs

Regression

Bias-Variance
Composition

- The regression function which minimises the expected squared loss, is given by the mean of the conditional distribution $p(t | \mathbf{x})$.



<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



Review

Linear Basis Function Models

Maximum Likelihood and Least Squares

Squared Loss

Regularized Least Squares

Multiple Outputs

Function for Regression

Bias-Variance Composition

- Analyse the expected loss

$$\mathbb{E}[L] = \iint \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt.$$

- Re

<https://eduassistpro.github.io>

$$= \{y(\mathbf{x}) - \mathbb{E}[t | \mathbf{x}]\}^2 + 2\{y(\mathbf{x}) - \mathbb{E}[t | \mathbf{x}]\} \{\mathbb{E}[t | \mathbf{x}] - t\}$$

- Claim

$$\iint \{y(\mathbf{x}) - \mathbb{E}[t | \mathbf{x}]\} \{\mathbb{E}[t | \mathbf{x}] - t\} p(\mathbf{x}, t) d\mathbf{x} dt = 0.$$

Add WeChat edu_assist_pro



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Function for
estimation

Bias-Variance
decomposition

- Claim

$$\iint \{f(\mathbf{x}) - \mathbb{E}[t | \mathbf{x}]\} \cdot \{\mathbb{E}[t | \mathbf{x}] - t\} p(\mathbf{x}, t) d\mathbf{x} dt = 0$$

- See
dep

<https://eduassistpro.github.io>

- Calculate the integral over t

$$\begin{aligned} \int \{\mathbb{E}[t | \mathbf{x}] - t\} p(\mathbf{x}, t) dt &= \mathbb{E}[t | \mathbf{x}] p(\mathbf{x}) - p(\mathbf{x}) \int t dt \\ &= \mathbb{E}[t | \mathbf{x}] p(\mathbf{x}) - p(\mathbf{x}) \mathbb{E}[t | \mathbf{x}] \\ &= 0 \end{aligned}$$

Add WeChat edu_assist_pro



Assignment Project Exam Help

- The expected loss is now

\mathbb{E}

- Min already).
- Second term represents the intrinsic variability of target data (can be regarded as noise). Independent choice $y(\mathbf{x})$, can not be reduced by learning a be

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Function for
Regression

Bias-Variance
Composition

The Bias-Variance Decomposition (1)

Statistical Machine
Learning

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



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Prediction

Bias-Variance
Decomposition

- Consider again squared loss for which the optimal prediction is given by the conditional expectation $\mathbb{E}(y|\mathbf{x})$.

- Since \mathcal{D} is a finite sample from the unknown joint distribution $p(\mathbf{x}, y)$, the learned function $f_{\mathcal{D}}$ depends on the sample \mathcal{D} .
- Notate the dependency of the learned function by $y(\mathbf{x}; \mathcal{D})$.
- Evaluate performance of algorithm by taking the expectation $\mathbb{E}_{\mathcal{D}} [L]$ over all data sets \mathcal{D} .

The Bias-Variance Decomposition (2)

Statistical Machine
Learning

© 2020

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



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Prediction

Bias-Variance
Decomposition

- Taking the expectation over data sets \mathcal{D} , using Eqn 1, and interchanging the order of expectations for the first term.

2

<https://eduassistpro.github.io>

- Again, add and subtract the expectation

$$\mathbb{E}_{\mathcal{D}} \left[\left(y(\mathbf{x}; \mathcal{D}) - \hat{y}(\mathbf{x}) \right)^2 \right] = \mathbb{E}_{\mathcal{D}} \left[\left(y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}} [y(\mathbf{x}; \mathcal{D})] + \mathbb{E}_{\mathcal{D}} [y(\mathbf{x}; \mathcal{D})] - \hat{y}(\mathbf{x}) \right)^2 \right]$$

and show that the mixed term vanishes under the expectation $\mathbb{E}_{\mathcal{D}}$.

The Bias-Variance Decomposition (3)



- Expected loss $\mathbb{E}_{\mathcal{D}} [L]$ over all data sets \mathcal{D}

expected loss = (bias)² + variance + noise.

where

<https://eduassistpro.github.io>

$$\text{noise} = \iint \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x}$$

- (bias)² : How accurate is a model across different data sets? (How much does the average prediction over all data sets differ from the desired regression function ?)
- variance : How sensitive is the model to small changes in the training set? (How much do solutions for individual data sets vary around their average ?)

The Bias-Variance Decomposition

Statistical Machine Learning

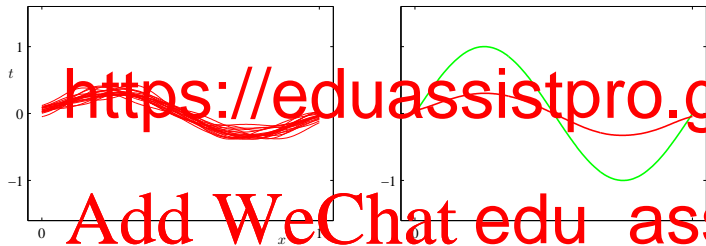
© 2020

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



Simple models have low variance and high bias.

Assignment Project Exam Help



Left: Result of fitting the model to 100 data s

Right: Average of the 100 fits in red, the sinusoidal function from where the data were created in green.

Review

Linear Basis Function Models

Maximum Likelihood and Least Squares

Squared Loss

Regularized Least Squares

Multiple Outputs

Selection for

Bias-Variance

Decomposition

The Bias-Variance Decomposition

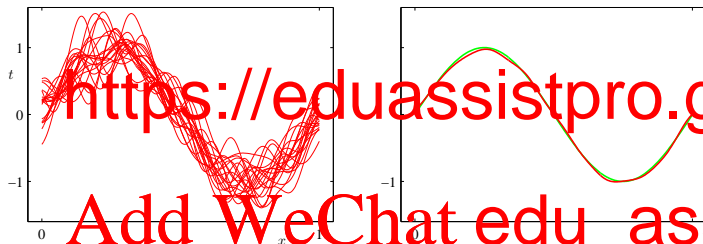
Statistical Machine
Learning

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



Complex models have high variance and low bias

Assignment Project Exam Help



Left: Result of fitting the model to 100 data s

Right: Average of the 100 fits in red, the sinusoidal function from where the data were created in green.

Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

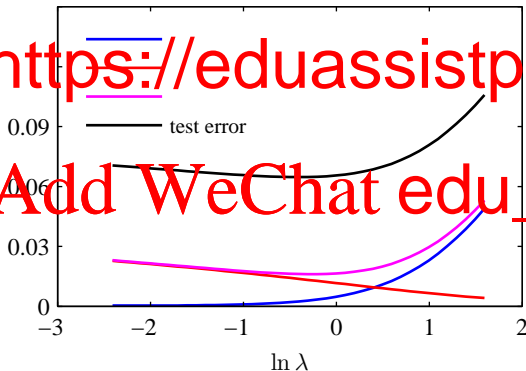
Multiple Outputs

Selection for
Bias-Variance

Decomposition

The Bias-Variance Decomposition

- Dependence of bias and variance on the model complexity
- Squared bias, variance, their sum, and test data
- The minimum for $(\text{bias})^2 + \text{variance}$ occurs close to the value that gives the minimum error





Assignment Project Exam Help

- You may have encountered *unbiased estimators*
- Why guarantee zero bias? To quote the pioneer of Bayesian inference, Edwin Jaynes, from his book *Pro*

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Prediction

Bias-variance
decomposition

The Bias-Variance Decomposition

Statistical Machine
Learning

© 2020

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



Review

Linear Basis Function
Models

Maximum Likelihood and
Least Squares

Squared Loss

Regularized Least
Squares

Multiple Outputs

Selection for
Prediction

Bias-variance
decomposition

- Tradeoff between bias and variance
 - simple models have low variance and high bias
 - complex models have high variance and low bias

- The
mo

- Exp

<https://eduassistpro.github.io>

$$\text{expected loss} = (\text{bias})^2 + \text{varia}$$

- The noise comes from the data, and can not be removed from the expected loss.
- To analyse the bias-variance decomposition : many data sets needed, which are not always available.

Add WeChat [edu_assist_pro](#)