



# Statistical Machine Learning

Christian Walder

Machine Learning Research Group  
CSIRO Data61

and

College of Engineering and Computer Science  
The Australian National University

Canberra  
Semester One, 2020.

## Outlines

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

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



## Part XXIII

# *Discussion and Summary*

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*

# Flavour of this course

- Formalise intuitions about problems
- Use language of mathematics to express models
- Geometry, vectors, linear algebra for reasoning
- Probabilistic models to capture uncertainty
- Calculus to identify good parameters
- Probabilistic inference
- Design and analysis of algorithms
- Numerical algorithms in python
- Understand the choices when designing machine learning methods



*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- Frequentist vs. Bayes approach
- Conditional Probability
- Bayes Theorem
- Discrete vs. continuous random variables
- Distributions (Gaussian, Bernoulli, Binomial)
- Multivariate Distributions
- Change of Variables
- Conjugate Priors

(Chapter 2)

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- Vector Space
- Matrix-Vector Multiplication = Linear Combination
- Projection
- Positive (Semi)-definite Matrix
- Rank, Determinant, Trace, Inverse
- Eigenvectors, Eigenvalues
- Eigenvector Decomposition
- Singular Value Decomposition
- Gradient Calculation

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- Gradient descent
- Stochastic (On-line) Gradient Descent
- Global vs. Local extremum
- Lagrange Multipliers
- Quadratic programming (quadratic objective function, linear constraints)
- Dynamic programming (e.g. Hidden Markov Model)

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- Joint Probability factorises.
- Conditional Independence
- Independence Structure or "the absence of edges"
- Directed, Undirected and Factor Graphs
- Bayesian Network, Blocked Path and  $d$ -separation
- Markov Random Field, (maximal) Cliques
- Factor Graphs are Bipartite Graphs

(Chapter 8)

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*

# What is Machine Learning?



## *Definition (Mitchell, 1998)*

A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*





- Examples
- General Setup
- Inductive Bias
- Restricted Hypothesis Space
- Importance of understanding the restrictions and whether they are appropriate
- Do not train on the test set

(Chapter 1)

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- Estimate best predictor = training = learning

Given data  $(x_1, y_1), \dots, (x_n, y_n)$ , find a predictor  $f_{\mathbf{w}}(\cdot)$ .

- 1 Identify the type of input  $x$  and output  $y$  data
- 2 Propose a mathematical model for  $f_{\mathbf{w}}$
- 3 Design an objective function or likelihood
- 4 Calculate the optimal parameter ( $\mathbf{w}$ )
- 5 Model uncertainty using the Bayesian approach
- 6 Implement and compute (the algorithm in python)
- 7 Interpret and diagnose results

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- Regression
- Classification: binary and multiclass
- Clustering and density estimation
- Sequence prediction
- Dimensionality reduction

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



Given the input space  $\mathbf{V}$ , input data  $\mathbf{x} \in \mathbf{V}$ , and a set of classes  $\mathcal{C} = \{\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_K\}$ .

- Discriminant Function  $f(\mathbf{x})$

$$f : \mathbf{V} \rightarrow \mathcal{C}$$

- Discriminant Model  $p(\mathcal{C}_k | \mathbf{x})$ , then use decision theory
- Generative Model  $p(\mathbf{x}, \mathcal{C}_k)$ , then use decision theory

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- Maximum Likelihood (ML)

$$\theta^* = \arg \max_{\theta} p(\mathcal{D} | \theta)$$

- Maximum a Posteriori (MAP)

$$\theta^* = \arg \max_{\theta} p(\theta | \mathcal{D}) \propto p(\mathcal{D} | \theta) p(\theta)$$

- Bayesian

$$p(\theta | \mathcal{D}^{(n)}) \propto p(\mathbf{x}^{(n)} | \theta) p(\theta | \mathcal{D}^{(n-1)})$$

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- Parametric Methods : Learn the model parameter from the training data, then discard training data.
- Nonparametric methods: Use training data for prediction
  - Histogram method
  - $k$ -nearest neighbours
  - Parzen probability density model: set of function centered on the data
- Kernel methods: Use linear combination of functions evaluated at the training data.

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- Bayes' Theorem

$$p(\boldsymbol{\theta} | \mathcal{D}) = \frac{p(\mathcal{D} | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathcal{D})}$$

- Regularized Risk

$$\sum_{n=1}^N \ell(\boldsymbol{\theta}, \mathcal{D}_n) + \frac{\lambda}{2} \|\boldsymbol{\theta}\|^2$$

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- General Regression Setup
- Closed-form solution
- Maximum Likelihood and Least Squares
- Geometry of Least Squares
- Sequential Learning (on-line)
- Choice of basis function
- Regularisation
- Powerful with nonlinear feature mappings
- Bias-Variance Decomposition

(Chapter 3.1, 3.2, 3.3)

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*





- Closed-form solution
- Predictive Distribution

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

- Conjugate Prior
- Limitations of Linear Basis Function Models
- Curse of dimensionality

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- General Classification Setup
- Input space versus Feature space
- Binary and Multiclass Labels
- Maximum Likelihood solution

$$\theta^* = \arg \max_{\theta} p(\mathbf{X}, \mathbf{t} | \theta)$$

- Naive Bayes : all features conditioned on the class are independent of each other

(Chapter 4)

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- Smooth logistic sigmoid acting on a linear feature vector
- Compare to perceptron
- Error as negative log likelihood (**cross-entropy** error)
- Gradient of error is target deviation times basis function (linear)
- Laplace approximation

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- Neural Networks
- Multilayer Perceptron with differentiable activation function
- The basis functions can now adapt to the data.
- Weight space symmetries.
- Error Backpropagation.
- Regularisation in Neural Networks.

(Chapter 5)

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- Maximise Variance
- Find the eigenvectors of the covariance corresponding to the largest eigenvalues.
- PCA and Compression
- Data Standardisation
- Data Whitening

(Chapter 12.1)

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- For feedforward neural networks, multiple layers can be advantageous
- Multiple PCA layers is equivalent to one single PCA
- Want to minimise reconstruction error, add nonlinear hidden layer
- Undercomplete autoencoder - lossy compression
- Pre-training of supervised learning (with unlabelled data)
- Denoising autoencoder
- Overcomplete autoencoder - sparse representations

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- Joint Probability over observed variables does not longer factorise.
- Introduce discrete latent variables to model complex marginal distributions over the observed variables by simpler distributions over observed and latent variables.
- $K$ -means clustering
- Data compression
- Mixture of Bernoulli
- Mixture of Gaussians

$$\begin{aligned} p(\mathbf{x}) &= \sum_{\mathbf{z}} p(\mathbf{z}) p(\mathbf{x} | \mathbf{z}) = \sum_{\mathbf{z}} \prod_{k=1}^K \pi_k^{z_k} \prod_{k=1}^K \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k} \\ &= \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \end{aligned}$$

(Chapter 9.1, 9.2)

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*

# Expectation Maximisation (EM)



- Evaluate the responsibilities, then maximise the parameters.
- **E step**: Find  $p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}})$ .
- **M step**: Find  $\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}})$  where

$$Q(\theta, \theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z} | \theta)$$

- Kullback-Leibler Divergence

(Chapter 9.3, 9.4)

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*





- Inner Product  $\rightarrow$  Kernel
- Kernels are a kind of similarity measure
- Sparse Kernel Machines
- Support Vector Machines
- Lagrange multipliers
- How do we get to the relevant data points?
- Overlapping class distribution
- Output are decisions, not posterior probabilities.

(Chapter 6.1, 6.2 and 7.1)

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

***Kernels***

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- Sum Rule

$$p(X) = \sum_Y p(X, Y)$$

- Four possible options to do  $\sum_Y$ 
  - Brute force - exhaustive
  - Sampling
  - Sum product
  - Variational Inference

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- Central task: Find  $p(\mathbf{Z} | \mathbf{X}, \theta)$ .
- Stochastic Approximation
- Sampling from the uniform distribution.
- Sampling from standard distributions via the inversion of the cumulative distribution function
- Rejection Sampling
- Adaptive Rejection Sampling
- Importance Sampling : calculate the expectation of function  $f(z)$  under distribution  $p(z)$  using some simpler  $q(z)$ .
- Markov Chain Monte Carlo
- Metropolis Hastings Algorithm
- Gibbs Sampling

(Chapter 11)

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- Independence structure - Local computations
- Sum-Product Algorithm
- Message passing
- Distributive Law

(Chapter 8.4)

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- Deterministic Approximation
- Assume an analytical approximation of posterior
- E.g. Assume posterior factorizes
- Convert integration problem into optimization problem over functions
- Use exponential family form, optimize with KL divergence

(Chapter 10)

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- Stationary vs. Nonstationary Sequential Distributions
- Markov Model of order  $M = 0, 1, \dots$
- State Space Model using latent variables
- Hidden Markov Model (HMM): Latent variables are discrete.
- Homogeneous HMM
- Left-to-right HMM
- Viterbi algorithm

(Chapter 13.1, 13.2)

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



- **The Language**

Vectors, geometry, linear algebra, calculus, probability, graphical models

- **The Tool**

Construct model, take the gradient, set it to zero, solve

- **Supervised learning**

Regression, classification and sequence prediction

- **Unsupervised learning**

Clustering, density estimation and dimensionality reduction

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*

# Flavour of this course



- Formalise intuitions about problems
- Use language of mathematics to express models
- Geometry, vectors, linear algebra for reasoning
- Probabilistic models to capture uncertainty
- Calculus to identify good parameters
- Design and analysis of algorithms
- Numerical algorithms in python
- Understand the choices when designing machine learning methods

*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*



# What we did not cover (in detail)

- Other learning paradigms
  - Neural Networks
  - Evolutionary Methods (e.g. Genetic Algorithms)
  - Frequent item mining
  - Expert Systems / Rule based learning
- Theory
  - Information Theory
  - Convex Optimisation
  - Generalised Linear Models
  - Dynamical Systems
  - Reinforcement Learning
  - Artificial Intelligence
- Applications
  - Natural Language Processing
  - Computer Vision
  - Computational Social Science
  - Robotics



*The Language*

*Problem Setting*

*Linear Regression and  
Classification*

*Neural Networks*

*Non-Factorising  
Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

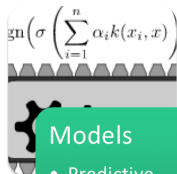
*Where to go from here?*

# What is Machine Learning?



## Data

- Vectors
- Structured



## Models

- Predictive functions
- Probabilistic Models



## Training

- Maximum Likelihood
- Gradient Descent

Shameless plug: [mml-book.com](http://mml-book.com)



*The Language*

*Problem Setting*

*Linear Regression and Classification*

*Neural Networks*

*Non-Factorising Distributions*

*Kernels*

*Sum Rule*

*Factorising Distributions*

*Sequential Data*

*Where to go from here?*