

# COMP6229 Machine Learning MSc

## Week 1: Introduction

### Mahesan Niranjan

Department of Electronics and Computer Science

#### Lecturers:

Mahesan Niranjan    Katayoun (Kate) Farrahi



Slides are prompts (for me); Notes are what you make, off the white-board and from textbooks during self study.

Autumn Semester 2017/18

Mahesan Niranjan (UoS)

COMP6229

Autumn Semester 2017/18

1 / 30

## Overview

- Logistics
- Motivation
  - Some examples from my research
- Review of Mathematical Foundations
  - Linear Algebra
  - Calculus
  - Probability Theory / Statistics
  - Principles of Optimization

Emphasis is on *foundations* of the subject (mathematical and algorithmic). We will not do formal mathematics here, instead we develop an understanding of the concepts and tools.

## Teaching:

- One lecture per week (some weeks two)
- One tutorial / catch up / help session per week
- Ten three-hour lab sessions
- One reading week
- One revision lecture before exam.

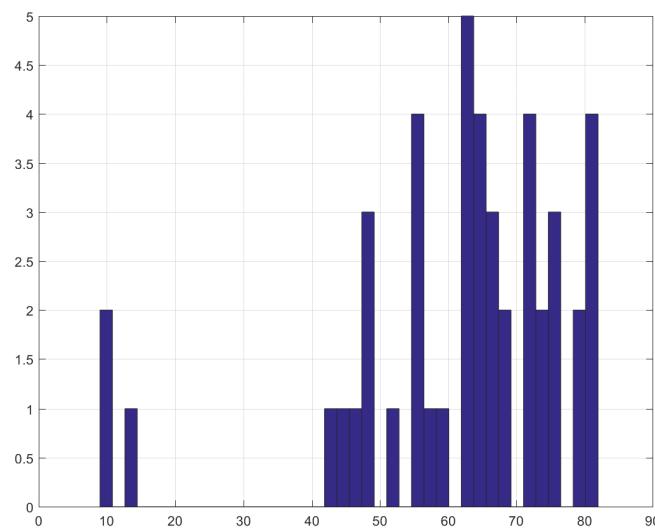
## Assessment:

20% Coursework (from ~W4)  
80% Semester end written exam

- MSc passmark 50%

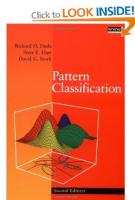
# Assessment

Distribution of marks, COMP6229 2015/16

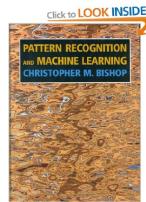


Difficult to fail this module, but please don't try!

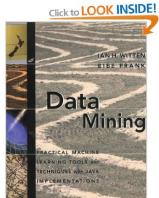
## Good Books



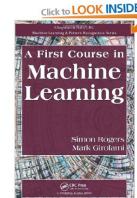
R.O.Duda, P.E.Hart & D.G.Stork  
Pattern Classification



C.M. Bishop  
Pattern Recognition and Machine Learning



I.H. Witten & E. Frank  
Data Mining



S. Rogers & M. Girolami  
A First Course in Machine Learning

*"There is nothing to be learnt from a professor, which is not to be met with in books"*  
- David Hume (1711-1776)

(WikiPedia: "Hume had little respect for the professors of his time [...] He did not graduate")

Machine Learning: Good employment prospects!

**Research Engineer - Software Developer - Machine Learning - London**  
Apply on your phone  
  
London ([/jobs/in-london](#)) | South East ([/jobs/in-south-east](#))  
From £70,000 to £120,000 per annum + bonus + benefits  
  
Miller Maxwell Ltd ([/jobs-at-miller-maxwell/jobs](#))  
Permanent  
Full-time  
Full-time  
 Miller M

at/miller  
Research Engineer - Software Developer - Machine Learning - C++ - Java - Python - London. Established, high profile technology driven **Proprietary Trading House** based in the City. No commercial finance experience is required for this exciting Research Engineer, Software Developer position.

## Trading House in City

The Institute of Cancer Research - Search Engine: Postdoctoral Training Fellow  
Ref:1594737

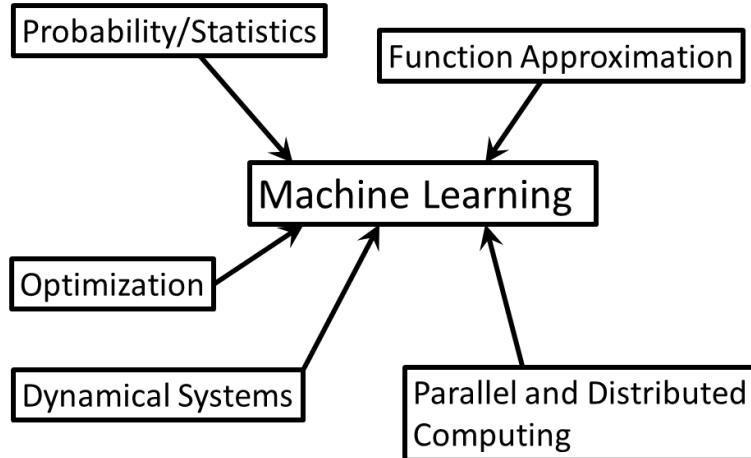
Postdoctoral Training Fellow – Image analysis (x2), - Ref:1594737

<a href="#">Click here to go back to search results</a>	
Closing Date of vacancy	30 Oct 2016
Division	Molecular Pathology
Team	Computational Pathology & Integrated Genomics
Type of Contract	Fixed Term
Length of Contract	3 years
Salary Range	<b>£29,960 - £42,820 p.a. inclusive (full salary scale)</b>
Work Location	Fusion (Surrey)
Hours per week	35

The Institute of Cancer Research, London, is one of the world's most influential cancer research institutes with an outstanding record of achievement dating back more than 100 years. We provided the first convincing evidence that DNA damage is the basic cause of cancer, laying the foundation for the now universally accepted idea that cancer is a genetic disease. Today, The Institute of Cancer Research (ICR) leads the world at isolating cancer-related genes and discovering

---

Standard disclaimers apply!



- Mathematical / Statistical side of Artificial Intelligence
- Machine Learning draws from many fields

## Machine Learning as Data-driven Modelling

Single-slide overview of the subject and challenging questions

Data  $\{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N \quad \{\mathbf{x}_n\}_{n=1}^N$

Function Approximator  $\mathbf{y} = f(\mathbf{x}, \boldsymbol{\theta}) + v$

Parameter Estimation  $E_0 = \sum_{n=1}^N \{\|\mathbf{y}_n - f(\mathbf{x}_n; \boldsymbol{\theta})\|\}^2$

Prediction  $\hat{\mathbf{y}}_{N+1} = f(\mathbf{x}_{N+1}, \hat{\boldsymbol{\theta}})$

Regularization  $E_1 = \sum_{n=1}^N \{\|\mathbf{y}_n - f(\mathbf{x}_n)\|\}^2 + g(\|\boldsymbol{\theta}\|)$

Modelling Uncertainty  $p(\boldsymbol{\theta} | \{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N)$

Probabilistic Inference  $\mathbf{E}[g(\boldsymbol{\theta})] = \int g(\boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta} = \frac{1}{N_s} \sum_{n=1}^{N_s} g(\boldsymbol{\theta}^{(n)})$

Sequential Estimation  $\boldsymbol{\theta}(n-1|n-1) \rightarrow \boldsymbol{\theta}(n|n-1) \rightarrow \boldsymbol{\theta}(n|n)$   
Kalman & Particle Filters; Reinforcement Learning

# Machine Learning

Many Interesting Problems (to me)

- Visual Scene Recognition
  - Machine Translation
  - Computational Biology
  - Computational Finance
  - Recommender Systems
  - Physiological Signal Modelling
- 
- “Big Data: ” Buzzword causing even more excitement!
  - Make accurate predictions and make money!
  - Make statements about the problem domain and become famous!

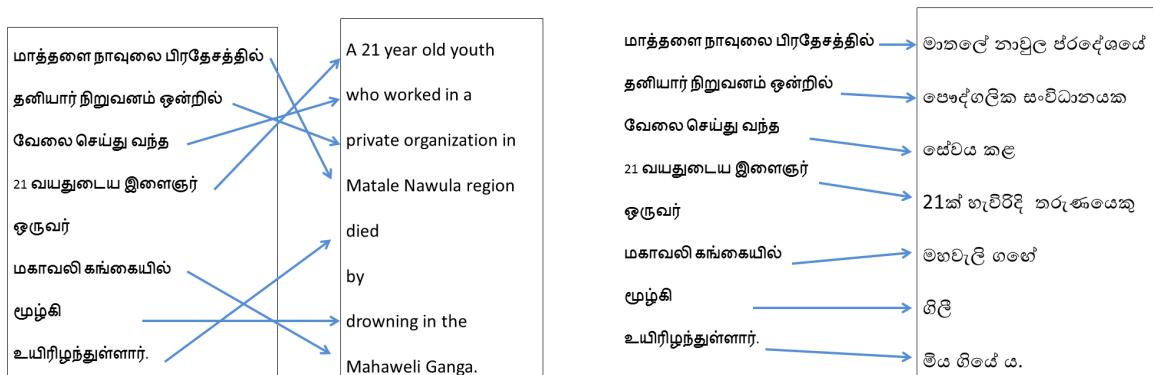
ECS: Advanced courses building on the foundations you will learn here:

- Advanced Machine Learning
- Computational Biology
- Computational Finance

## Examples from my research

Example 1: Machine Translation

- Phrases move due to grammatical differences.
- Variability due to context of phrase.



- Not data rich (electronically available parallel corpora);
- Solution from active learning.

# Examples from my research

## Example 2: Computational Finance

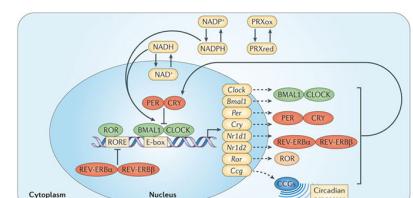
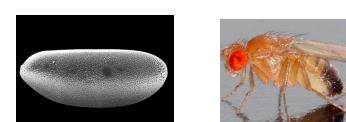
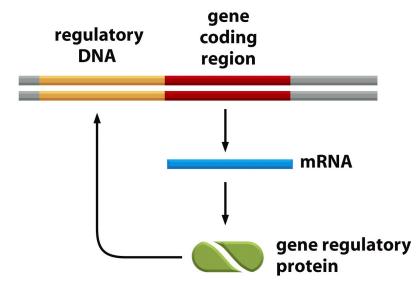
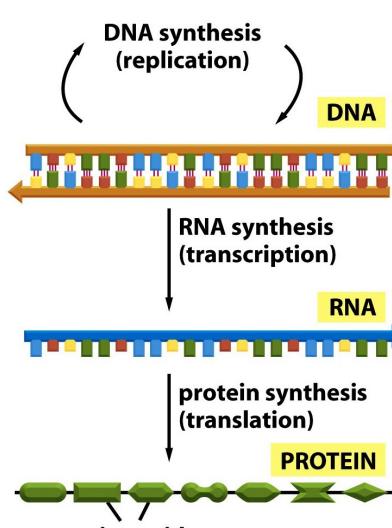
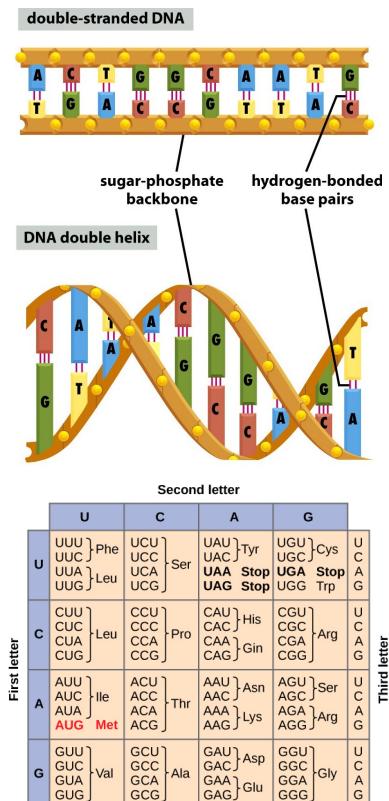
### Constructing Sparse Portfolios

A. Takeda, M. Niranjan, J. Gotoh & Y. Kawahara (2013) "Simultaneous pursuit of out-of-sample performance and sparsity in index tracking portfolios", Computational Management Science 10(1): 21-49.

See White Board

## Molecular Biology

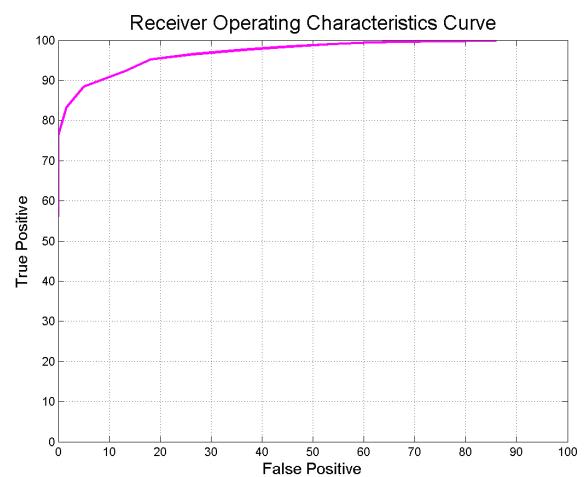
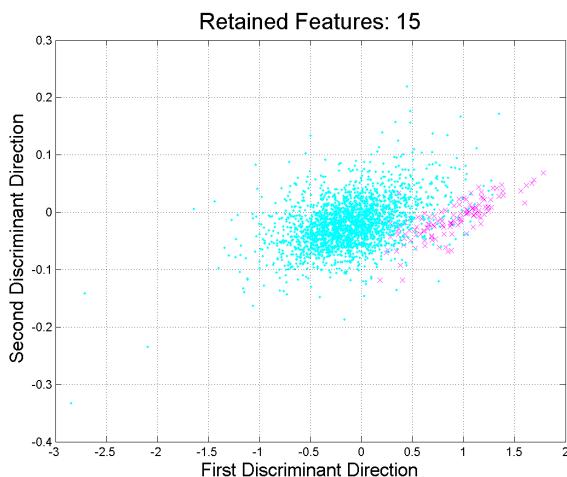
(Figures from: Alberts *et al.* Molecular Biology of the Cell)



# Examples from my research

## Example 3: Classifying Gene Function

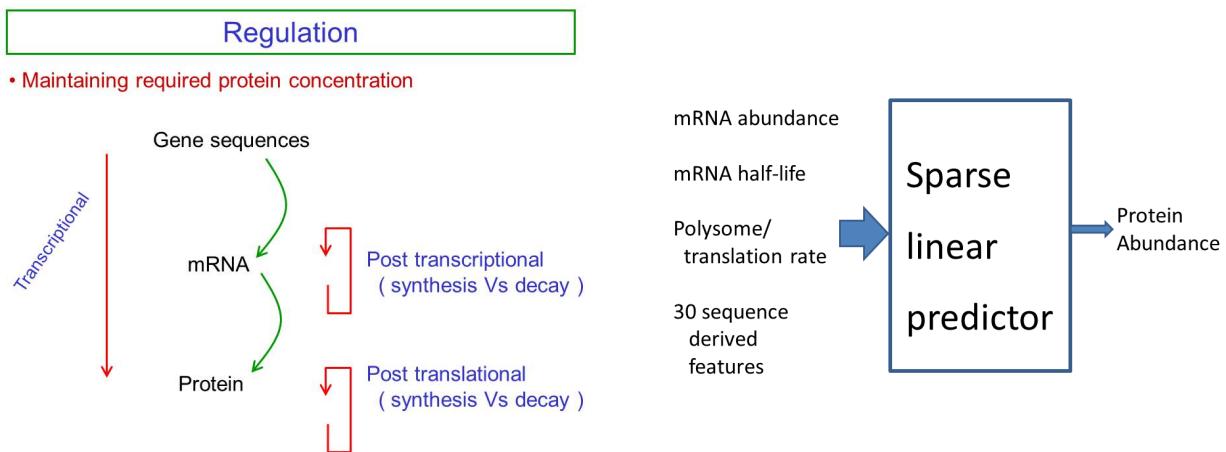
- 2000 yeast genes
- Observed (simultaneously) under 78 conditions
- Some have a specific function; others not



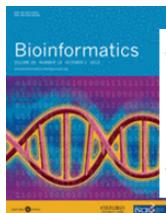
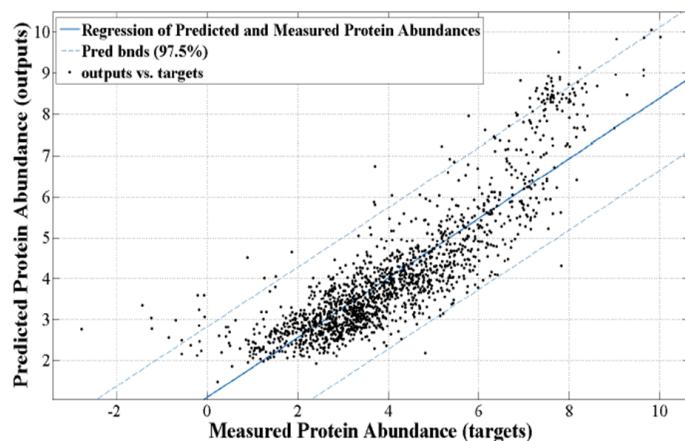
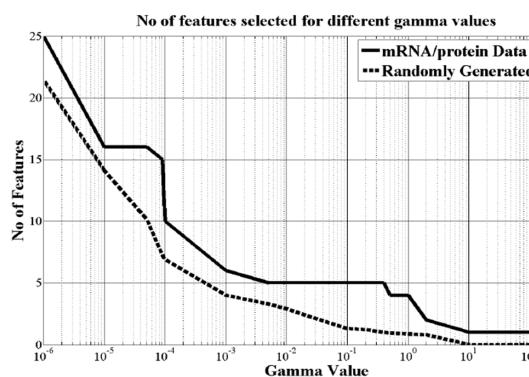
See MATLAB Demo

## Example 4: Regulation of Protein Concentrations [Yawwani]

Gunawardana]



- Set up a predictor of protein concentration
- Sparse model selects relevant features
- Outliers  $\Rightarrow$  post-translationally regulated proteins



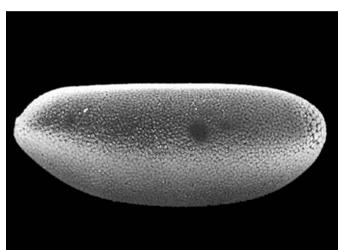
Systems biology

Bridging the gap between transcriptome and proteome measurements identifies post-translationally regulated genes

Yawwani Gunawardana and Mahesan Niranjan\*

[doi:10.1093/bioinformatics/btt537](https://doi.org/10.1093/bioinformatics/btt537)

## Example 5: Morphogen Propagation in Development

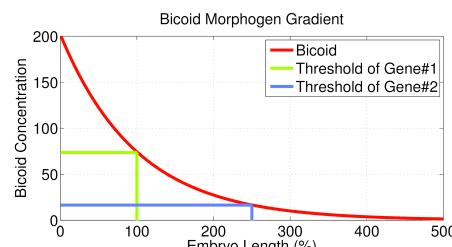


A. Turing

C. Nüsslein-Volhard

$$\frac{\partial}{\partial t} M(x, t) = D \frac{\partial^2}{\partial x^2} M(x, t) - \tau_p^{-1} M(x, t) + S(x, t)$$

B. Houchmandzadeh et al. (2007), *Nature*



**Establishment of developmental precision and proportions in the early *Drosophila* embryo**

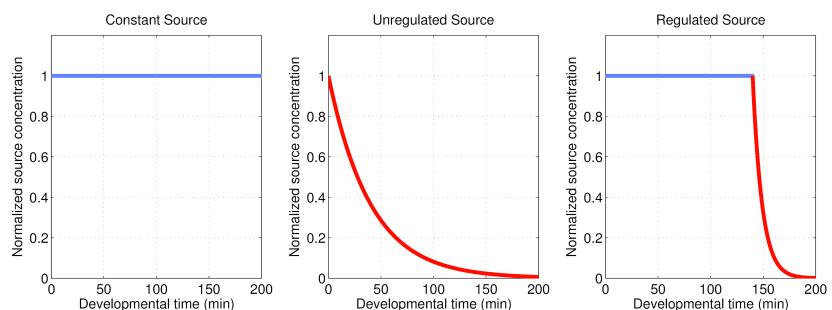
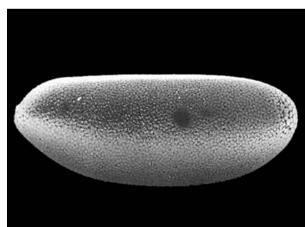
Bahram Houchmandzadeh\*†, Eric Wieschaus\* & Stanislas Leibler\*‡§

\* Howard Hughes Medical Institute, Department of Molecular Biology, Princeton University, Princeton, New Jersey 08544, USA

† CNRS, Laboratoire de Spectrométrie Physique, BP87, 38402, St-Martin D'Hères Cedex, France

# Is Maternal mRNA Stability Regulated?

Wei Liu



OPEN ACCESS Freely available online

PLOS one

## The Role of Regulated mRNA Stability in Establishing Bicoid Morphogen Gradient in *Drosophila* Embryonic Development

Wei Liu\*, Mahesan Niranjan

School of Electronics and Computer Science, University of Southampton, Southampton, United Kingdom

### Abstract

The Bicoid morphogen is amongst the earliest triggers of differential spatial pattern of gene expression and subsequent cell

BIOINFORMATICS

Vol. 00 no. 00 2005  
Pages 1–6

### Gaussian process modelling for *bicoid* mRNA regulation in spatio-temporal Bicoid profile

Wei Liu and Mahesan Niranjan\*

School of Electronics and Computer Science, University of Southampton, Southampton, UK

Mahesan Niranjan (UoS)

COMP6229

Autumn Semester 2017/18

17 / 30

## Example 6: Systems Level Modelling

[Xin Liu]

$$\begin{aligned}\dot{D}_t &= K_d \frac{S_t}{1 + \frac{K_s D_t}{1 + K_u U_f}} - \alpha_d D_t \\ \dot{S}_t &= \eta(t) - \alpha_0 S_t - \alpha_s \frac{\frac{K_s D_t}{1 + K_u U_f}}{1 + \frac{K_s D_t}{1 + K_u U_f}} S_t \\ \dot{U}_f &= K(t)[P_t - U_f] - [K(t) + K_{fold}] D_t\end{aligned}$$



Systems biology

Vol. 28 no. 11 2012, pages 1501–1507

### State and parameter estimation of the heat shock response system using Kalman and particle filters

Xin Liu and Mahesan Niranjan\*

School of Electronics and Computer Science, University of Southampton, Southampton SO17 1BJ, UK  
Associate Editor: Trey Ideker

Mahesan Niranjan (UoS)

COMP6229

Autumn Semester 2017/18

18 / 30

- Linear Algebra

- Calculus

- Optimization

- Probabilities

- This is not a course on any of the above!

- We need tools from these topics.

- Quickly review what we need today, and will return to each topic as and when we need them (in just about enough depth) to understand machine learning.

## Linear Algebra: Vectors and Matrices

- Vectors and matrices as collections of numbers

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1d} \\ a_{21} & a_{22} & \dots & a_{2d} \\ \vdots & \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nd} \end{bmatrix}$$

- Operations on collections on numbers

- Scalar product

$$\mathbf{w} \bullet \mathbf{x} = \sum_{i=1}^n w_i x_i$$

- With useful geometric insights

- Angle between vectors in  $n$  dimensional space

$$\mathbf{w} \bullet \mathbf{x} = |\mathbf{w}| |\mathbf{x}| \cos(\theta)$$

# Vectors

- Linear independence

... set of  $p$  vectors  $\mathbf{x}_j, j = 1, \dots, p$

$$\sum_{i=j}^p \alpha_j \mathbf{x}_j = \mathbf{0}$$

... only solution is all  $\alpha_j = 0$

... no vector in the set can be expressed as a linear combination of the others.

- Scalar product as projection: projection of vector  $\mathbf{x}$  on a direction specified by vector  $\mathbf{u}$

$$\frac{\mathbf{x} \bullet \mathbf{u}}{|\mathbf{u}|^2} \mathbf{u}$$

... we will also write this as

$$\frac{\mathbf{x}^T \mathbf{u}}{|\mathbf{u}|^2} \mathbf{u}$$

# Matrices

- Simple operations e.g. addition:  $[A + B]_{ij} = [A]_{ij} + [B]_{ij}$ ;  
transpose:  $[A]_{ij}^T = [A]_{ji}$ ; multiplication by a scalar:  $[\alpha A]_{ij} = \alpha [A]_{ij}$
- Matrix multiplication:

$$[AB]_{ij} = \sum_{k=1}^n [A]_{ik} [B]_{kj}$$

- $(AB)^T = B^T A^T$
- Square: number of rows = number of columns
- Symmetric:  $A^T = A$
- Identity matrix:  $I$  diagonal elements 1, off diagonals 0.
- Determinant:  $\det \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} = a_{11} a_{22} - a_{21} a_{12}$
- Trace:  $\text{trace}(A) = \sum_{i=1}^n a_{ii}$

# Linear transformation

- $\mathbf{y} = A\mathbf{x}$
- Rotation:  $R = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$   
 $R\mathbf{x}$  rotates  $\mathbf{x}$  by angle  $\theta$  radians.  
Magnitude of  $\mathbf{x}$  does not change.
- A special relationship between a square matrix  $A$  and vector  $\mathbf{x}$

$$A\mathbf{x} = \lambda\mathbf{x}$$

Magnitude scales, but no rotation... have you come across this?

Eigenvalues, eigenvectors

Found by

$$\det(A - \lambda I) = 0$$

- **Homework:** Look up if the following are true and how they are proved.

- $\det(A) = \prod_{i=1}^n \lambda_i$
- $\text{trace}(A) = \sum_{i=1}^n \lambda_i$

- Real symmetric matrix  $A = U D U^T$   
Columns of  $U$  orthogonal.

- More advanced (very powerful) topic: Singular value decomposition (SVD)

## Rapid Review of Foundations II: Calculus

- Function  $y = f(\mathbf{x})$ 
  - Derivative  $\frac{dy}{dx}$  is gradient/slope;
  - Integral  $\int_{x=a}^{x=b} f(x)dx$  is area under the curve.
- Function of several variables  $y = f(x_1, x_2, \dots, x_p)$ 
  - Partial derivatives  $\frac{\partial f}{\partial x_i}$ : Differentiate with respect to  $x_i$  pretending all other variables remain constant.
  - Gradient vector

$$\nabla f = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_p} \end{pmatrix}$$

**Homework:** Consider  $f = \mathbf{x}^t A \mathbf{x}$ ,  $A = A^T$ ;  $\nabla f = 2A\mathbf{x}$ . Using scalars in two dimensions, i.e.  $\mathbf{x} = [x_1 \ x_2]^T$  and  $A$  to contain elements  $\begin{bmatrix} a_{11} & a_{12} \\ a_{12} & a_{22} \end{bmatrix}$ , verify the claim. Writing out the algebra helps in learning!

# Rapid Review of Foundations III: Optimization

- Unconstrained optimization:  $\min f(\mathbf{x})$

- Constrained optimization:

$$\begin{aligned} & \min f(\mathbf{x}) \\ & \text{subject to } g_i(\mathbf{x}) \leq b_i, \quad i = 1, 2, \dots, m \end{aligned}$$

- Gradient and Hessian:

$$\nabla \mathbf{f} = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_p} \end{pmatrix} \quad \mathbf{H} = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_p} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_p} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_p \partial x_1} & \frac{\partial^2 f}{\partial x_p \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_p^2} \end{bmatrix}$$

## Optimizations (cont'd)

- Example: Gradient descent algorithm

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \eta \nabla \mathbf{f}$$

- Newton's Method

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{H}^{-1} \nabla \mathbf{f}$$

- Example: Lagrange Multipliers

$$\begin{aligned} & \min f(\mathbf{x}) \\ & \text{subject to } g_i(\mathbf{x}) \leq b_i, \quad i = 1, 2, \dots, m \end{aligned}$$

$$F(\mathbf{x}, \lambda) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i [b_i - g_i(\mathbf{x})]$$

- We will use various optimization algorithms in this module (later in the coursework).
- Advanced **Homework:** Search for CVX Disciplined Convex Programming and have a rough read.

## Rapid Review of Foundations IV: Probabilities

- Discrete probabilities  $P[X]$
- Continuous densities  $p(x)$
- Joint  $P[X, Y]$ ; Marginal  $P[X]$ ; Conditional  $P[X|Y]$

$$P[Y|X] = \frac{P[X|Y] P[Y]}{P[X]}$$
$$P[X] = \sum_Y P[X|Y] P[Y]$$

$$P[X, Y] = P[X|Y] P[Y]$$
$$P[X] = \sum_y P[X, Y]$$
$$= \sum_Y P[X|Y] P[Y]$$

## Gaussian Densities: Univariate and Multivariate

- Univariate Gaussian

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2} \frac{(x - m)^2}{\sigma^2}\right\}$$

What are properties we know? **Homework:** Draw sketches for different values of  $m$  and  $\sigma$ .

- Multivariate Gaussian

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{p/2}(\det \mathbf{C})^{1/2}} \exp\left\{-\frac{1}{2} (\mathbf{x} - \mathbf{m})^t \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m})\right\}$$

Mean  $\mathbf{m}$  is a vector

Covariance,  $\mathbf{C}$ , matrix: symmetric, positive semi definite!

**Homework:** Draw sketches for different values of  $\mathbf{m}$  and  $\mathbf{C}$

$$\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{C}), \mathbf{y} = \mathbf{Ax} \implies \mathbf{y} \sim \mathcal{N}(\mathbf{Am}, \mathbf{A}\mathbf{C}\mathbf{A}^T)$$

---

Univariate Mean       $\hat{m} = \frac{1}{N} \sum_{n=1}^N x_n$

Univariate Covariance     $\hat{\sigma} = \frac{1}{N} \sum_{n=1}^N (x_n - \hat{m})^2$

Multivariate Mean       $\hat{\mathbf{m}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$

Covariance Matrix       $\hat{\mathbf{C}} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \hat{\mathbf{m}})(\mathbf{x}_n - \hat{\mathbf{m}})^T$

---

- These are known as maximum likelihood estimates (see later).
- **Homework:** Have you noticed there are two buttons in a calculator for estimating standard deviation, denoted  $\sigma_n$  and  $\sigma_{n-1}$ ? Find out why.

## What Next?

Weeks 2 and 3: Pattern Classification

- Classifying based on  $P[\omega_j | \mathbf{x}]$
- Optimal classifier for simple distributions
- Linear classifier – when is it optimal?
- Distance based classifiers
  - Nearest Neighbour classifier
  - Mahalanobis distance
- Linear discriminant analysis
  - Fisher LDA
- Classifier Performance
  - Receiver Operating Characteristics (ROC) Curve
- Perceptron learning rule and convergence

See sketches on whiteboard – these illustrations are important