# Advanced probabilistic methods Lecture 1

Pekka Marttinen

Aalto University

January, 2019

Pekka Marttinen (Aalto University)

- Practical matters
  - Structure, workload, grading
  - Exercise format
  - Student feedback from 2018
- Course overview
- Basic probability calculus (Barber, Ch. 1)
- Basic graph concepts

Pekka Marttinen (Aalto University)

<sup>&</sup>lt;sup>1</sup>These slides build upon the book *Bayesian Reasoning and Machine Learning* and the associated teaching materials. The book and the demos can be downloaded from *www.cs.ucl.ac.uk/staff/D.Barber/brml.* 

#### Structure

- Lectures  $10 \times 2$  hours
- $\bullet~$  Exercise sessions  $(1+9)\times 2$  hours
- See Timetable in myCourses/Materials

#### Grading based on

- Exam: 70% of the total weight
- Exercises: 30%
- Preliminary boundaries: 1:50%, 2:60%, 3:70%, 4:80%, 5:90%
- Minimum required: 1/2 of the exam points, 1/3 of exercise points





• Note: *Bayesian Reasoning and Machine Learning* is freely available for download at *www.cs.ucl.ac.uk/staff/D.Barber/brml* 

- Lectures:  $10 \times 2h$
- Preparation for lectures, reading the book ( $\sim$  200 pages): 9  $\times$  4h
- Exercise sessions: 9 × 2h
- Doing the exercises:  $9 \times 4.5h$
- Exercise 0, self-study, introduction to Python: 5.5h
- Preparing for the exam: 11h
- Exam: 4h
- Total 135h. As credits 135/27 = 5cr.

• Exercises must be returned to MyCourses by the deadline

- A single PDF
- Grading of the exercises (2p→done, almost correct; 1p→done, but something clearly missing/incorrect; 0p→not done or completely incorrect)
- Exercises are graded by the TAs, not corrected → Always make sure afterwards you know the correct answer, by attending the exercise sessions or going through the model solutions.
- Exercise session format
  - help for getting started with next week's exercises
  - possibility to ask about next week's exercises or previous week's solutions
  - two assistants present

#### Relation to other courses



#### About prerequisites

 'Bayesian Data Analysis should be listed as a prerequisite' →See previous slide.

#### About exam

- 'Exam was difficult'
  - $\rightarrow$ Exam duration increased from 3 to 4 hours.

 $\rightarrow \textsc{Overall}$  level of performance was (and will be) taken into account in grading.

 $\rightarrow$  Questions in the exam will be similar to the exercise questions. Best way to prepare is to do the exercises.

 ${\rightarrow}\mathsf{Separate}$  recap lecture added to the end of the course.

- About difficulty in general
  - 'At least to me there was steep curve in the topic difficulty starting from lecture 5. Compress lectures 1-4 and spend more time on the "new stuff"

 $\rightarrow \mathsf{Take}$  a full advantage of the Exercise sessions (also and especially on the 2nd half)

 $\rightarrow$ Ask clarifications on the lectures

- About exercises
  - 'Participating in the exercise sessions was extremely useful and the course assistants were great'
  - 'The exercises really helped me to understand the concepts related to this course, and I think that they played a huge role in my learning'
  - 'TAs from the exercise session were very helpful and stayed often (a lot) longer!'

#### Student feedback from 2018, overview

#### • Overall assessment



#### • I will benefit from things learnt on the course



Pekka Marttinen (Aalto University)

January, 2019 10 / 41

- The goal of **probabilistic modeling** is to answer a question about the data:
  - Classify the samples into groups
  - Create prediction for future observations
  - Select between competing hypotheses
  - Estimate a parameter, such as the mean, of the population

۰...



# Probabilistic modeling overview (2/2)

- Probabilistic modeling in a nutshell
  - Select a model
  - Infer the parameters of the model (train/fit the model)
  - Use the fitted model to answer the question of interest
- Usually several models are considered, requiring model selection.
- For example:  $f_n(x) = a_0 + a_1 x + a_2 x^2 + \ldots + a_n x^n$



# Course contents (1/2)

• Ingredients of probabilistic modeling

- **Models**: Bayesian networks, Sparse Bayesian linear regression, Gaussian mixture models, latent linear models
- Methods for inference: maximum likelihood, maximum a posteriori (MAP), Laplace approximation, expectation maximization (EM), Variational Bayes (VB), Stochastic variational inference (SVI)
- Ways to select between models



# Course contents (2/2)

• A brief introduction to probabilistic programming using Tensorflow and Edward<sup>2</sup>.





# Google

<sup>2</sup>Tensorflow Probability is another very recent tool that incorporates also the functionality of Edward, see: https://www.tensorflow.org/probability/. $\equiv$  + < $\equiv$  +

Pekka Marttinen (Aalto University)

Advanced probabilistic methods

January, 2019 14 / 41

#### Role of probabilistic machine learning today

• Keynote at NIPS in December 2017 by Yee Whye Teh



https://www.youtube.com/watch?v=9saauSBgmcQ

Image: Image:

- Marginalization
- Independence
- Conditional distribution
- Conditional independence
- Continuous random variables

(To recap these, see Additional Reading in myCourses/Materials)

- Random variables: X, Y, Z, ...
- Values these random variables can take: x, y, z, ...
- Probability
  - The following notations are used interchangeably

$$p(X = x) = p_X(x) = p(x)$$

• All are interpreted as the probability that variable X is in state x

- Domain
  - dom(X) denotes all possible states for variable X.
- Distribution of a variable X consists of
  - its domain dom(X)
  - and full specification of probability values  $p_X(x)$ , for all possible  $x \in dom(X)$
- Normalization
  - The summation over all the states

$$\sum_{x \in dom(X)} p(X = x) = 1$$

• The sum can be written as:  $\sum_{x} p(x) = 1$ 

#### Example - probability table

В	М	Κ	p(b, m, k)
1	1	1	0.012
1	1	0	0.108
1	0	1	0.288
1	0	0	0.192
0	1	1	0.016
0	1	0	0.064
0	0	1	0.096
0	0	0	0.224

- The probability table lists the probabilities of all possible combinations of the random variables.
- The *joint* distribution of *B*, *M* and *K*

• For example

$$p_{B,M,K}(1,1,0) = p(B=1, M=1, K=0)$$
  
= 0.108

• Modified from Example 1.3 "Inspector Clouseau"

M = 'Maid is the murderer'

B = 'Butler is the murderer'

K = 'Knife is the murder weapon'

• Given a joint dist  $p_{X,Y}(x, y)$ , the marginal dist of X is defined by

$$p_X(x) = \sum_{y} p_{X,Y}(x,y)$$

• More generally,

$$p(x_1,...,x_{i-1},x_{i+1},...,x_n) = \sum_{x_i} p(x_1,...,x_n)$$

- В K p(b, m, k)Μ 0.012 1 1 1 1 1 0 0.1081 0 1 0.288 1 0 0 0.192 0 1 0.016 1 0 1 0 0.064 0 0 1 0.096 0 0 0 0.224
- What is the marginal distribution of *B* and *M*?
- We need to compute  $p_{B,M}(b, m)$ , for all possible *b* and *m*.

# Example - marginalization (2/2)

Use

$$p_{B,M}(b,m) = \sum_{k=0}^{1} p_{B,M,K}(b,m,k)$$

• For example:

$$p_{B,M}(0,0) = p_{B,M,K}(0,0,0) + p_{B,M,K}(0,0,1)$$
  
= 0.096 + 0.224 = 0.32

• Doing this for all *B*, *M* combinations, we get the marginal probability table

В	Μ	p(b, m)
1	1	0.12
1	0	0.48
0	1	0.08
0	0	0.32

• Random variables X and Y are independent if

$$p_{X,Y}(x,y) = p_X(x)p_Y(y)$$

for all x and y.

- Intuitively, this means that knowing the value of X does not provide any information about the value of Y.
- Notation: X ⊥ Y
- More generally:  $\mathcal{A} = \{A_1, \dots, A_k\}$  and  $\mathcal{B} = \{B_1, \dots, B_l\}$  are independent if

$$p_{A_1,...,A_k,B_1,...,B_l}(a_1,...,a_k,b_1,...,b_l) = p_{A_1,...,A_k}(a_1,...,a_k)p_{B_1,...,B_l}(b_1,...,b_l)$$

В	Μ	p(b,m)
1	1	0.12
1	0	0.48
0	1	0.08
0	0	0.32

• Are *B* and *M* independent?

#### Marginal distributions

В	p(b)		Μ	p(m)
1	0.6	and	1	0.2
0	0.4		0	0.8

Direct computation gives

В	Μ	p(b)p(m)	p(b,m)
1	1	0.12	0.12
1	0	0.48	0.48
0	1	0.08	0.08
0	0	0.32	0.32

• Hence, B and M are (marginally) independent

- D='number of people drowned', A='amount of ice-cream sold'
  - Are D and A independent?
  - Are D and A causally dependent?

#### Conditional distribution

$$p_{X|Y}(x|y) = \frac{p_{X,Y}(x,y)}{p_Y(y)}$$

specifies the probability of each possible value x of X given that we have observed variable Y in state y.

#### Example - Conditional distribution

• For example:

$$p(K = 1|B = 1, M = 1) = \frac{p(B = 1, M = 1, K = 1)}{p(B = 1, M = 1)} = 0.1$$
  
 $p(K = 0|B = 1, M = 1) = 0.9$ 

• All conditional probabilities in the last column

В	Μ	Κ	p(b, m, k)	p(k b,m)
1	1	1	0.012	0.1
1	1	0	0.108	0.9
1	0	1	0.288	0.6
1	0	0	0.192	0.4
0	1	1	0.016	0.2
0	1	0	0.064	0.8
0	0	1	0.096	0.3
0	0	0	0.224	0.7

 X ⊥ Y | Z denotes that variables X and Y are conditionally independent of each other, given the state of variable Z. This is formally defined by condition

$$p_{X,Y|Z}(x,y|z) = p_{X|Z}(x|z)p_{Y|Z}(y|z)$$

for all states x, y, z of variables X, Y, Z.

• Intuitively, this means that if we know the value of Z, knowing in addition the value of Y does not provide any information about the value of X. Indeed, provided p(y, z) > 0, we have

$$X \perp Y | Z \Longrightarrow p_{X|Y,Z}(x|y,z) = p_{X|Z}(x|z)$$

$$X \perp Y | Z \Longrightarrow p_{X|Y,Z}(x|y,z) = p_{X|Z}(x|z)$$

Proof

$$p(x|y, z) = \frac{p(x, y, z)}{p(y, z)} = \frac{p(x, y|z)p(z)}{p(z)p(y|z)}$$
$$= \frac{p(x|z)p(y|z)p(z)}{p(z)p(y|z)} = p(x|z)$$

• The general chain rule of probability

$$p(x, y, z) = p(x|y, z)p(y|z)p(z),$$

follows from iterative use of the definition of conditional probability.

- В Κ p(b, m, k)Μ 1 1 1 0.012 1 0.108 1 0 1 0 1 0.288 1 0 0 0.1920 1 0.016 1 0 1 0 0.064 0 0 1 0.096 0 0 0 0.224
- Are *M* and *B* conditionally independent, given *K*?
- We need to compare
  - $p_{M|K}(m|k)p_{B|K}(b|k)$
  - $p_{B,M|K}(b,m|k)$

for all m, b, k.

# Example - Conditional independence (2/3)

• For example,

$$p(B = 1, M = 1 | K = 1) = \frac{p(B = 1, M = 1, K = 1)}{p(K = 1)}$$
$$= \frac{0.012}{0.012 + 0.288 + 0.016 + 0.096} \approx 0.0291$$

• Similarly,

$$p(M = 1 | K = 1) = \frac{p(M = 1, K = 1)}{p(K = 1)}$$
$$= \frac{0.012 + 0.016}{0.012 + 0.288 + 0.016 + 0.096} \approx 0.0508$$

and

$$p(B=1|K=1)=\ldots\approx 0.7110$$

Pekka Marttinen (Aalto University)

э

В  $K \quad p(b,m|k) \quad p(b|k) \quad p(m|k) \quad p(b|k)p(m|k)$ Μ 0.029 0.711 0.051 0.036 1 1 1 0 1 1 . . . . . . . . . . . . 1 0 1 0 0 1 1 1 0 0 1 0 1 0 0 0 0 0

• Because 0.029  $\neq$  0.036, it follows that *B* and *M* are not conditionally independent given *K*.

# Intuition for independence and conditional independence $\left( 1/2 \right)$

- Let  $X_1, X_2, \ldots, X_n$  denote the cumulative sum of n dice throws, such that  $dom(X_1) = \{1, \ldots, 6\}$ ,  $dom(X_2) = \{2, \ldots, 12\}$ , etc.
  - Is  $X_{n+1}$  independent of  $X_{n-1}$ ?
  - Is  $X_{n+1}$  conditionally independent of  $X_{n-1}$  given  $X_n$ ?
- X='Location of an airplane now', Y='Location of the plane 15s ago', Z='Location 15s from now'
  - Is Y independent of Z?
  - Is Y conditionally independent of Z given X?

# Intuition for independence and conditional independence $\left(2/2\right)$

- S='sunshine', D='number of people drowned', A='amount of ice-cream sold'
  - Are *D* and *A* independent?
  - Are D and A conditionally independent given S?
- A='The alarm is on', B=There is a burglar in the house", T='A truck passes the house'
  - Suppose that the alarm can be triggered either by a burglar or by a passing truck
  - Are *B* and *T* independent?
  - Are B and T conditionally independent given A

# Continuous random variables (1/3)

• Probability density function (pdf) for a continuous variable X,  $f_X()$ 

$$\int_{x \in \mathcal{R}} f_X(x) dx = 1$$
$$p(X \in [a, b]) = \int_{x=a}^{b} f_X(x) dx$$

Cumulative distribution function (cdf)

$$F_X(x) = p(X \le x) = \int_{t=-\infty}^{x} f_X(t) dt$$

σ<sup>2</sup>=0.2, - $\sigma^2 = 1.0,$ 

 $\sigma^2 = 5.0$ 





- Concepts presented can be generalized to continuous random variables
- Marginalization
  - Discrete:  $p_X(x) = \sum_y p_{X,Y}(x,y)$
  - Continuous:  $f_X(x) = \int_y f_{X,Y}(x,y) dy$
- Expected value
  - Discrete:  $E(X) = \sum_{x} x p_X(x)$
  - Continuous:  $E(X) = \int_X x f_X(x) dx$

Conditional distribution

$$f_{Y|X}(y|x) = \frac{f_{X,Y}(x,y)}{f_X(x)}$$

• (conditional) independence:  $X \perp Y | Z$ , if

$$f_{X,Y|Z}(x,y|z) = f_{X|Z}(x|z)f_{Y|Z}(y|z)$$

# Basic graph definitions



- A graph consists of **nodes** (vertices) and undirected of directed **edges** (links) between nodes.
- A path from X<sub>i</sub> to X<sub>j</sub> is a sequence of connected nodes starting at X<sub>i</sub> and ending at X<sub>j</sub>.



A Directed Acyclic Graph (DAG) is a directed graph without cycles
Parents, Children, Ancestors, Descendants,... (see Ch. 2)

- marginalization
- conditional distribution
- conditional/marginal independence
- probability density function, cumulative distribution function
- Basic graph concepts