		What is machine learning all about?	2(41)
Welcome to Ma	achine Learning 2014!! Thomas Schön Division of Systems and Control Department of Information Technology Uppsala University. Email: thomas.schon@it.uu.se, www: user.it.uu.se/~thosc112	"Machine learning is about learning, reasoning and acting based on data."	g
Machine Learning, Lecture 1 – Introduction T. Schön, 2014		Machine Learning, Lecture 1 – Introduction T. Schön, 2014	
Outline lecture 1	3(41)	Problem classes	4(41)
Outline lecture 1 1. Introduction and some r 2. Course administration 3. Probability distributions 1. Exponential family 2. Properties of the mu 3. Maximum Likelihood 4. Bayesian modeling 5. Robust statistics ("he 6. Mixture of Gaussians	motivating examples and some basic ideas ltivariate Gaussian I (ML) estimation eavy tails") s	 Problem classes Supervised learning. The training data consists of both input and output (target) data. Classification: Discrete output variables. Regression: Continuous output variables. Unsupervised learning. The training data consists of input data only. Clustering: Discover groups of similar examples in data. Reinforcement learning. Finding suitable actions (control signals) in a given situation in order to maximize a reward. Close to control theory. This course is focused on supervised learning. 	4(41) It

Example 1 – autonomous helicopter aerobatics 5(41)	Example 2 – handwritten digit classification 6(41)
Field of the states and Andrew Y. Ng. Autonomous helicopter aerobatics through apprenticeship learning. International Journal of Robotics Research (IJRR), 29(13):1608-1639, November 2010. Interesting (very recent) developments on learning controllers are available here, Marc P. Deisenroth, Dieter Fox, and Carl E. Rasmussen Gaussian processes for data-efficient learning in robotics and control. IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMM), 2014. Machine Learning, Lecture 1 – Introduction T. Schön, 2014	 Input data: 16 × 16 grayscale images. Task: classify each input image as accurately as possible. This data set will be used throughout the course. Solutions and their performance are summarized on yann.lecun.com/ exdb/mnist/ Data set available from http://statweb.stanford.edu/~tibs/ElemStatLearn/
Example 3 – BNP for dynamical systems 7/41)	Example 4 – animal detection and tracking (I/II)
<text><list-item><list-item><list-item><list-item></list-item></list-item></list-item></list-item></text>	<text><section-header><section-header><text><text><text><text><text><text><text><text><text><text><text><text><text><text><text><text><text><text><text><text><text></text></text></text></text></text></text></text></text></text></text></text></text></text></text></text></text></text></text></text></text></text></section-header></section-header></text>

Machine Learning, Lecture 1 – Introduction T. Schön, 2014 Machine Learning, Lecture 1 – Introduction T. Schön, 2014

Example 4 – animal detection and trac	cking (II/II) 9(41)	Field of machine learning 10(41)
Were automatikka skydd mot påköming år under utveckling. Snart ska det kunna ugpstakka skora djur ja konsektion. Not ska skydd mot påköming år under utveckling. Snart ska det kunna ugpstakka skora djur ja konsektion.	 Learning detectors for animals. Sensor fusion between radar and infrared camera. 	 Top 3 conferences on general machine learning Neural Information Processing Systems (NIPS) International Conference on Machine Learning (ICML) Uncertainty in Artificial Intelligence (UAI) and Inter. Conf. on Artificial Intelligence and Statistics (AISTATS) Top 3 journals on general machine learning Journal of Machine Learning Research (JMLR) IEEE Trans. Pattern Analysis and Machine Intelligence (PAMI) IEEE Trans. on Neural Networks (TNN) and Neural computation (NECO)
Machine Learning, Lecture 1 – Introduction T. Schön, 2014		Machine Learning, Lecture 1 – Introduction T. Schön, 2014
Course administration	11(41)	Course administration – projects (3 hp) 12(41)
 Course administration Lecturer and examiner: Thomas Schön 11 lectures (do not cover everything) We will try to provide examples of active rest the lectures (especially connections to "our" Suggested exercises are provided for each Written exam, 3 days (72 hours). Code of husual All course information, including lecture mat from the course home page http://user.it.uu.se/~thoscl12/ML/ 	search throughout ' areas) lecture onor applies as terial is available	 2(41) Ocurse administration – projects (3 hp) Voluntary and must be based on a data set. Project ideas: discuss with me for ideas or even better, make up your own!! Form teams (1-3 students/project). Project time line: Mar. 24 Project proposal presentation May 2 Final reports are due See course home page for details. Note that the deadline for NIPS is in the beginning of June.

Project example from the 2011 edition

Detection and classification of cars in video images

Task: Train a detector/classifier, which can be used to detect, track and eventually classify different vehicles in the video recordings.

13(41)

14(41)

Calibration of a magnetometer in combination with inertial sensors

Project example from the 2013 edition

Task: Compute the calibration parameters of the magnetometer.





Multivariate Gaussian (II/VI) 21(41)	Multivariate Gaussian (III/VI) 22(41)
Theorem (Conditioning) Let x be Gaussian distributed and partitioned $x = \begin{pmatrix} x_a^T & x_b^T \end{pmatrix}^T$, then the conditional density $p(x_a x_b)$ is given by $p(x_a x_b) = \mathcal{N}(x_a \mu_{a b}, \Sigma_{a b}),$ $\mu_{a b} = \mu_a + \Sigma_{ab} \Sigma_{bb}^{-1}(x_b - \mu_b),$ $\Sigma_{a b} = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba},$ which using the information (precision) matrix can be written, $\mu_{a b} = \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab}(x_b - \mu_b),$ $\Sigma_{a b} = \Lambda_{aa}^{-1}.$	Theorem (Marginalization)Let x be Gaussian distributed and partitioned $x = (x_a^T \ x_b^T)^T$, then the marginal density $p(x_a)$ is given by $p(x_a) = \mathcal{N}(x_a \mid \mu_a, \Sigma_{aa}).$
Machine Learning, Lecture 1 – Introduction T. Schön, 2014	Machine Learning, Lecture 1 – Introduction T. Schön, 2014
Multivariate Gaussian (IV/VI) 23(41) Theorem (Affine transformations) Assume that x_a , as well as x_b conditioned on x_a , are Gaussian distributed $p(x_a) = \mathcal{N}(x_a \mid \mu_a, \Sigma_a),$ $p(x_b \mid x_a) = \mathcal{N}(x_b \mid Mx_a + b, \Sigma_{b a}),$ where M is a matrix and b is a constant vector. The marginal density of x_b is then given by $p(x_b) = \mathcal{N}(x_b \mid \mu_b, \Sigma_b),$ $\mu_b = M\mu_a + b,$ $\Sigma_b = \Sigma_{b a} + M\Sigma_a M^T.$	Multivariate Gaussian (V/VI) Theorem (Affine transformations, cont.) The conditional density of x_a given x_b is $p(x_a \mid x_b) = \mathcal{N}(x_a \mid \mu_{a b}, \Sigma_{a b}),$ with $\mu_{a b} = \Sigma_{a b} \left(M^T \Sigma_{b a}^{-1}(x_b - b) + \Sigma_a^{-1} \mu_a \right)$ $= \mu_a + \Sigma_a M^T \Sigma_b^{-1}(x_b - b - M\mu_a),$ $\Sigma_{a b} = \left(\Sigma_a^{-1} + M^T \Sigma_{b a}^{-1} M \right)^{-1}$ $= \Sigma_a - \Sigma_a M^T \Sigma_b^{-1} M \Sigma_a.$

Multivariate Gaussian (VI/VI) 25(41)	Maximum Likelihood (ML) estimation 26(41)
	Maximum likelihood provides a systematic way of computing point estimates of the unknown parameters θ in a given model, by exploiting the information present in the measurements $\{x_n\}_{n=1}^N$.
Multivariate Gaussian's are important building blocks in more sophisticated models. For more details, proofs and an example where the Kalman filter is derived using the above theorems is provided here, http://user.it.uu.se/~thosc112/pubpdf/schon12011.pdf	 Computing ML estimates of the parameters in a model amounts to: 1. Model the obtained measurements x₁,, x_N as a realisation from the stochastic variables x₁,, x_N. 2. Decide on which model to use. 3. Assume that the stochastic variables x₁,, x_N are conditionally iid. In ML the parameters θ are chosen in such a way that the measurements {x_n}^N_{n=1} are as likely as possible, i.e.,
Machine Learning, Lecture 1 – Introduction T. Schön, 2014	Machine Learning, Lecture 1 – Introduction T. Schön, 2014
Bayesian modeling 27(41)	Conjugate priors – example 1 (I/II) 28(41)
Bayesian modeling27(41)The goal in Bayesian modeling is to compute the posterior $p(\theta \mid x_{1:N})$.Provided that it makes sense from a modeling point of view it is convenient to choose prior distributions rendering a computationally tractable posterior distribution.This leads to the so called conjugate priors (if the prior and the posterior have the same functional form, the prior is said to be a conjugate prior for the likelihood).Again, only make use of conjugate priors if this makes sense from a modeling point of view!	Conjugate priors – example 1 (I/II) Let $X = \{x_n\}_{n=1}^N$ be independent identically distributed (iid) observations of $x \sim \mathcal{N}(\mu, \sigma^2)$. Assume that the variance σ^2 is known. The likelihood is given by $p(X \mid \mu) = \prod_{n=1}^N p(x_n \mid \mu) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2\right)$ If we choose the prior as $p(\mu) = \mathcal{N}(\mu \mid \mu_0, \sigma_0^2)$, the posterior will also be Gaussian. Hence, this Gaussian prior is a conjugate prior for the likelihood.

Conjugate priors – example 1 (II/II)

29(41)

The resulting posterior is

$$p(\mu \mid X) = \mathcal{N}(\mu_B, \sigma_B^2),$$

where the parameters are given by

$$\begin{split} \mu_{B} &= \frac{\sigma^{2}}{N\sigma_{0}^{2} + \sigma^{2}}\mu_{0} + \frac{N\sigma_{0}^{2}}{N\sigma_{0}^{2} + \sigma^{2}}\mu_{\text{MLL}} \\ \frac{1}{\sigma_{R}^{2}} &= \frac{1}{\sigma_{0}^{2}} + \frac{N}{\sigma^{2}}. \end{split}$$

The ML estimate of the mean is

$$\mu_{\mathsf{ML}} = \frac{1}{N} \sum_{n=1}^{N} x_n$$

Machine Learning, Lecture 1 – Introduction T. Schön, 2014

Conjugate prior is just one of many possibilities! 31(4)

Note that using a conjugate prior is **just one** of the many possible choices for modeling the prior! If it makes sense, use it, since it leads to simple calculations.

Let's have a look at an example where we do not make use of the conjugate prior and end up in a useful and interesting result.

Linear regression models the relationship between a continuous target variable t and an (input) variable x according to

$$t_n = w_0 + w_1 x_{1,n} + w_2 x_{2,n} + \dots + w_D x_{D,n} + \epsilon_n$$

= $w^T \phi(x_n) + \epsilon_n$,

where
$$\phi(x_n) = \begin{pmatrix} 1 & x_{1,n} & \dots & x_{D,n} \end{pmatrix}^T$$
 and $n = 1, \dots, N$.

Likelihood	Model Parameters	Conjugate Prior
Normal (known mean)	Variance	Inverse-Gamma
Multivariate Normal	Precision	Wishart
(known mean)		
Multivariate Normal	Covariance	Inverse-Wishart
(known mean)		
Multivariate Normal	Mean and covariance	Normal-Inverse-
		Wishart
Multivariate Normal	Mean and precision	Normal-Wishart
Exponential	Rate	Gamma

Machine Learning, Lecture 1 – Introduction T. Schön, 2014

Conjugate prior is just one of many possibilities! 32(41)

Let $\epsilon_n \sim \mathcal{N}(0, \sigma^2)$, resulting in the following likelihood

$$p(t_n \mid w) = \mathcal{N}(t_n \mid w^T \phi(x_n), \sigma^2).$$

Let us now assume w_n to be independent and Laplacian distributed (i.e. not conjugate prior), $w_n \sim \mathcal{L}(0, 2\sigma^2/\lambda)$

Def. (Laplacian distribution) $\mathcal{L}(x \mid a, b) = \frac{1}{2b} \exp\left(-\frac{|x-a|}{b}\right)$.

The resulting MAP estimate is given by,

$$\widehat{w}^{\mathsf{MAP}} = \operatorname*{arg\,max}_{w} \sum_{n=1}^{N} (t_n - w^T \phi(x_n))^2 + \lambda \sum_{n=1}^{D} |w_n|$$

Known as the **LASSO** and it leads to sparse estimates.

Machine Learning, Lecture 1 – Introduction T. Schön, 2014



Outlier detection models 37(41)	Summary – robust statistics 38(41)
Model the data as if it comes from a mixture of two Gaussians, $p(x_i) = p(x_i k_i = 0)p(k_i = 0) + p(x_i k_i = 1)p(k_i = 1)$ $= \mathcal{N}(0, \sigma^2)p(k_i = 0) + \mathcal{N}(0, \alpha\sigma^2)p(k_i = 1).$ where $\alpha > 1$, $p(k_i = 0)$ is the probability that the sample is ok and $p(k_i = 1)$ is the probability that the sample is an outlier. Note the similarity between these two "robustifications": • The Student's t-distribution is an infinite mixture of Gaussians, where the mixing is controlled by the ν -parameter. • The outlier detection model above consists of a sum of two Gaussians.	 Do not use distributions with thin tails (non-robust) if there are outliers present. Use more realistic robust "heavy tailed" distribution such as the Student's t-distribution or simply a mixture of two Gaussians. A nice account on robustness in a computer vision context is available in Section 3.1 in B. Triggs, P. McLauchlan, R. Hartley, and A. Fitzgibbon. Bundle Adjustment - A Modern Synthesis. In: Vision algorithms: theory and practice. Lecture Notes in Computer Science, Vol 1883:152–177. Springer, Berlin, 2000. dx. doi.org/10.1007/3-540-44480-7_21
Machine Learning, Lecture 1 – Introduction T. Schön, 2014	Machine Learning, Lecture 1 – Introduction T. Schön, 2014
Example – range measurements with outliers 39(41)	Important message! 40(41)
We measure range (r), contaminated by a disturbance $d_n \ge 0$ and noise $e_n \sim \mathcal{N}(0, \sigma^2)$, $y_n = r + d_n + e_n$. Compute the MAP estimate of $\theta = \{r, d_1, \dots, d_N\}$ under an exponential prior on d_n , $p(d_n) = \begin{cases} \lambda \exp(-\lambda d_n), & d_n \ge 0, \\ 0, & d_n < 0. \end{cases}$ Resulting problem $\widehat{\theta}^{MAP} = \arg \max p(\theta \mid y_{1:N}) = \arg \min \sum N \frac{(y_n - r - d_n)^2}{\sigma^2} + \lambda \sum_{n=1}^N d_n$	Given the computational tools that we have today it can be rewarding to resist the Gaussian convenience!!
 θ θ n=1 0² n=1 For details, see Example 2.2. in the PhD thesis of Jeroen Hol. This principle is used for ultra-wideband positioning, incorporated into MotionGrid (www.xsens.com/en/general/motiongrid) from our partners Xsens (www.xsens.com). 	We will try to repeat and illustrate this message throughout the course using theory and examples.
Machine Learning, Lecture 1 – Introduction T. Schön, 2014	Machine Learning, Lecture 1 – Introduction T. Schön, 2014

A few concepts to summarize lecture 1

41(41)

Supervised learning: The data consists of both input and output signals (e.g., regressions and classification).

Unsupervised learning: The data consists of output signals only (e.g., clustering).

Reinforcement learning: Finding suitable actions (control signals) in a given situation in order to maximize a reward. (Very similar to control theory)

Conjugate prior: If the posterior distribution is in the same family as the prior distribution, the prior and posterior are *conjugate distributions* and the prior is called a conjugate prior for the likelihood.

Maximum likelihood: Choose the parameters such that the observations are as likely as possible.

Machine Learning, Lecture 1 – Introduction T. Schön, 2014