

Probabilistic & Unsupervised Learning

**Week 1: Introduction, Statistical Basics,
and a bit of Information Theory**

Maneesh Sahani

`maneesh@gatsby.ucl.ac.uk`

**Gatsby Computational Neuroscience Unit, and
MSc in Intelligent Systems, Dept Computer Science
University College London**

Term 1, Autumn 2008

What is Learning?



Finding **structure** (**regularities**, **associations**) in observations.

Predicting new observations.

What is Learning?

Ideas related to learning appear in many fields:

- **Scientific Method:** epistemology, verification, experimental design, ...
- **Statistics:** theory of learning, data mining, learning and inference from data, ...
- **Computer Science:** AI, computer vision, information retrieval, ...
- **Engineering:** signal processing, system identification, adaptive and optimal control, information theory, robotics, ...
- **Cognitive Science:** computational linguistics, philosophy of mind, ...
- **Economics:** decision theory, game theory, operational research ...
- **Psychology:** perception, movement control, reinforcement learning, mathematical psychology...
- **Computational Neuroscience:** neuronal networks, neural information processing...

Different fields, Convergent ideas

- The **same set of ideas and mathematical tools** have emerged in many of these fields, albeit with different emphases.
- **Machine learning** is an interdisciplinary field focusing on both the mathematical foundations and practical applications of systems that learn, reason and act.
- **The goal of this course:** to introduce basic concepts, models and algorithms in machine learning with particular emphasis on unsupervised learning.

Three Types of Learning

Imagine an organism or machine which experiences a series of sensory inputs:

$$x_1, x_2, x_3, x_4, \dots$$

Supervised learning: The machine is also given **desired outputs** y_1, y_2, \dots , and its goal is to learn to **produce the correct output** given a new input.

Unsupervised learning: The goal of the machine is to **build a model** of x that can be used for reasoning, decision making, predicting things, communicating etc.

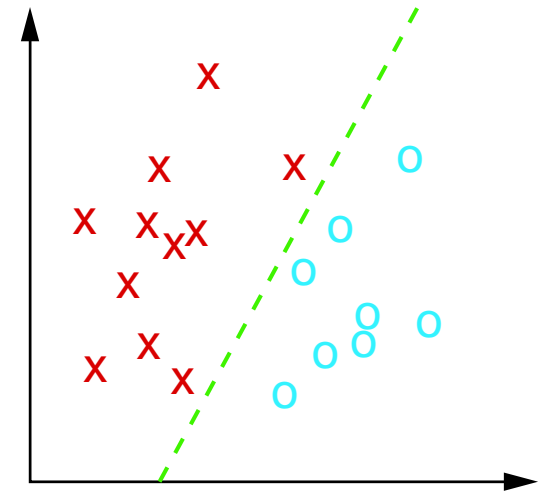
Reinforcement learning: The machine can also produce **actions** a_1, a_2, \dots which affect the state of the world, and receives **rewards (or punishments)** r_1, r_2, \dots . Its goal is to learn to act in a way that **maximises rewards** in the long term.

Goals of Supervised Learning

Two main examples:

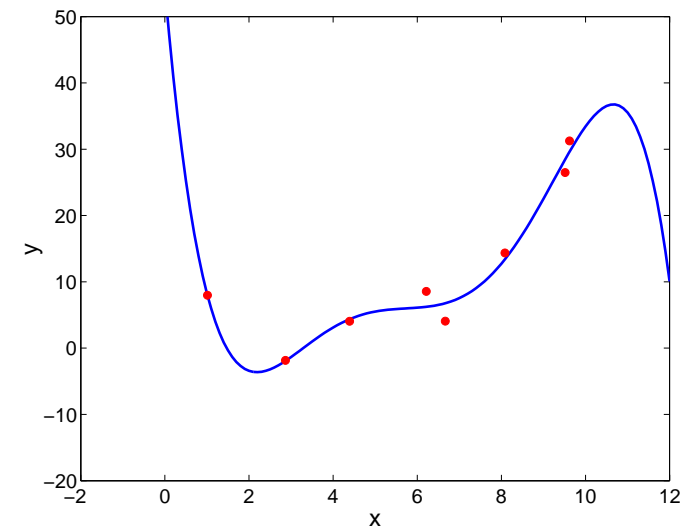
Classification:

The desired outputs y_i are discrete class labels.
The goal is to classify new inputs correctly
(i.e. to *generalize*).



Regression:

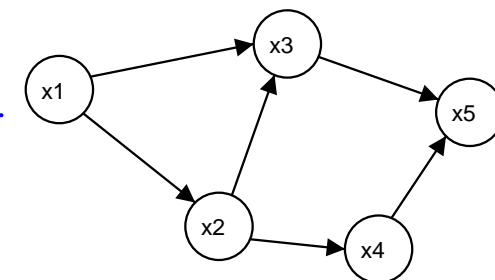
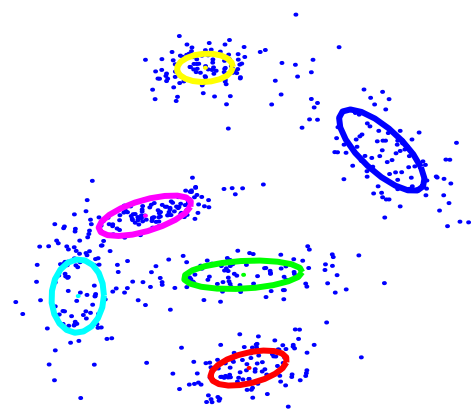
The desired outputs y_i are continuous valued.
The goal is to predict the output accurately for new inputs.



Goals of Unsupervised Learning

To build a model or find useful representations of the data, for example:

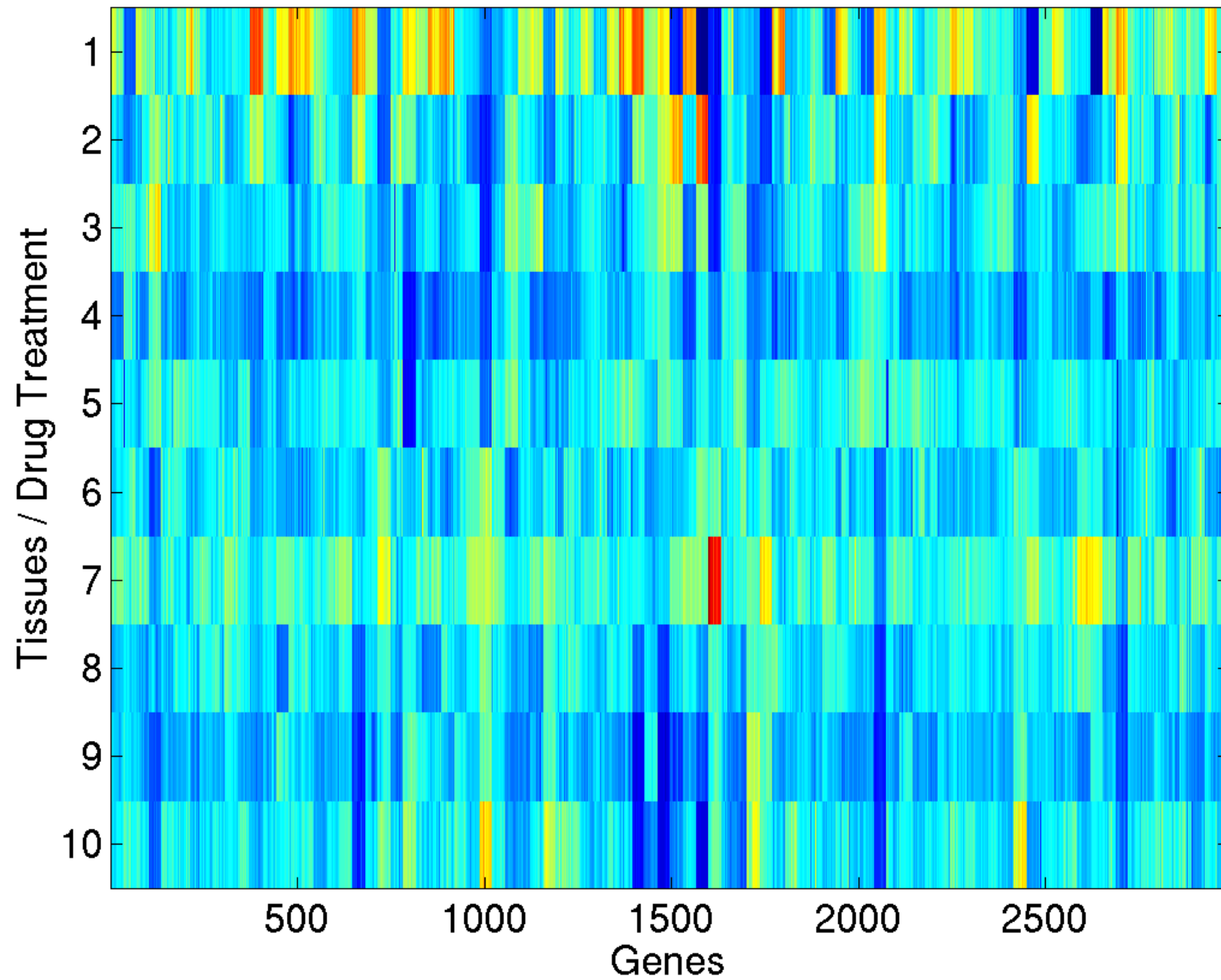
- finding clusters
- dimensionality reduction
- finding good explanations (hidden causes) of the data
- modeling the data density



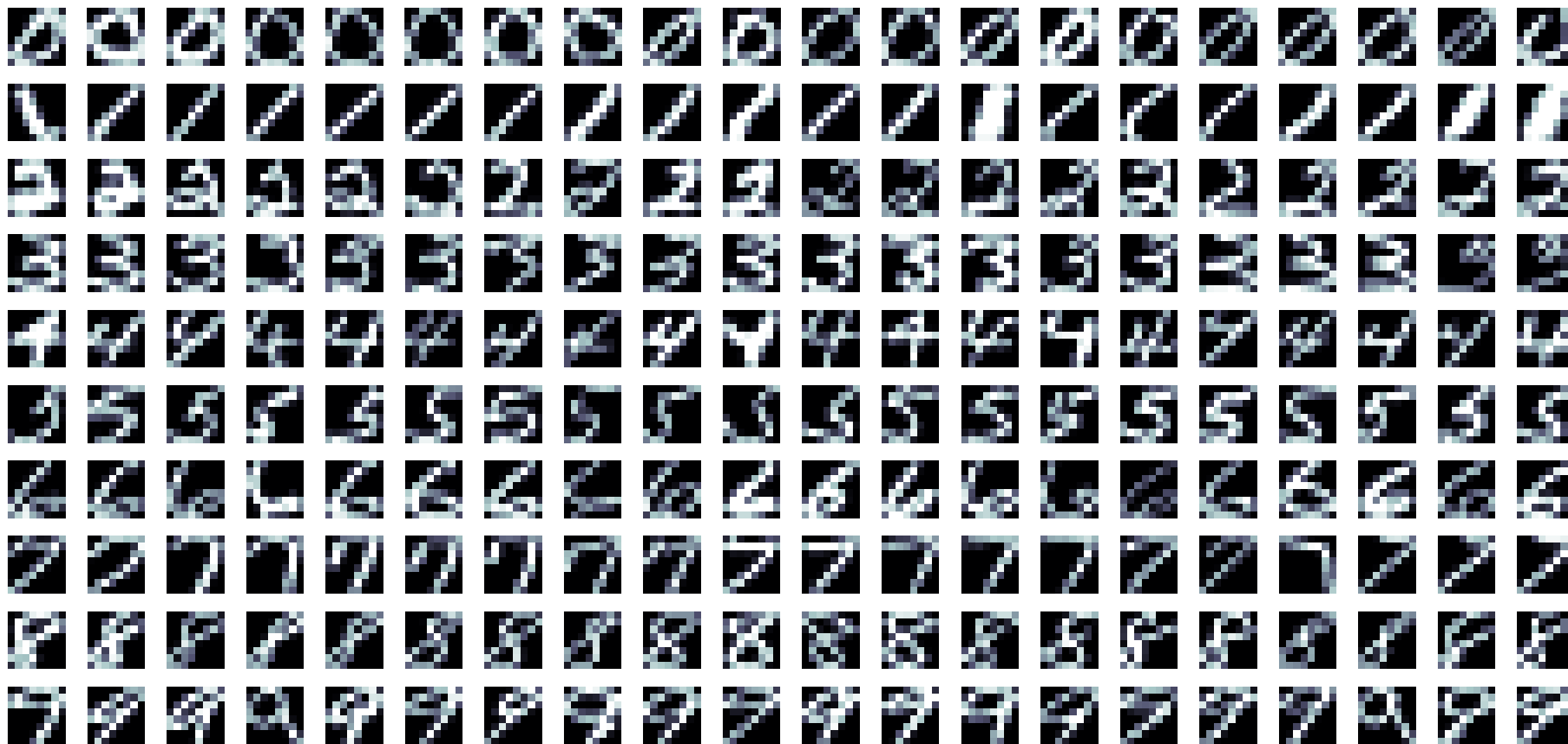
Uses of Unsupervised Learning

- structure discovery, science
- data compression
- outlier detection
- help classification
- make other learning tasks easier
- use as a theory of human learning and perception

Example data: Gene Expression



Example data: Handwritten Digits



Example data: Web Pages

Google Search: Unsupervised Learning http://www.google.com/search?q=Unsupervised+Learning&sourceid=fir...

 [Web](#) [Images](#) [Groups](#) [News](#) [Froogle](#) [more »](#)

Unsupervised Learning Search [Advanced Search](#) [Preferences](#)

Web Results 1 - 10 of about 150,000 for **Unsupervised Learning**. (0.27 seconds)

[Mixture modelling, Clustering, Intrinsic classification ...](#)
Mixture Modelling page. Welcome to David Dowe's clustering, mixture modelling and **unsupervised learning** page. Mixture modelling (or ...
www.csse.monash.edu.au/~did/mixture.modelling.page.html - 26k - 4 Oct 2004 - [Cached](#) - [Similar pages](#)

[ACL'99 Workshop -- Unsupervised Learning in Natural Language ...](#)
PROGRAM. ACL'99 Workshop **Unsupervised Learning** in Natural Language Processing. University of Maryland June 21, 1999. Endorsed by SIGNLL ...
www.ai.sri.com/~kebler/unsup-acl-99.html - 5k - [Cached](#) - [Similar pages](#)

[Unsupervised learning and Clustering](#)
cgm.cs.mcgill.ca/~soss/cs644/projects/wijhe/ - 1k - [Cached](#) - [Similar pages](#)

[NIPS*98 Workshop - Integrating Supervised and Unsupervised ...](#)
NIPS*98 Workshop "Integrating Supervised and **Unsupervised Learning**" Friday, December 4, 1998. ... 4:45-5:30. Theories of **Unsupervised Learning** and Missing Values. ...
www-2.cs.cmu.edu/~mccallum/supunsup/ - 7k - [Cached](#) - [Similar pages](#)

[NIPS Tutorial 1999](#)
Probabilistic Models for **Unsupervised Learning** Tutorial presented at the 1999 NIPS Conference by Zoubin Ghahramani and Sam Roweis. ...
www.gatsby.ucl.ac.uk/~zoubin/NIPStutorial.html - 4k - [Cached](#) - [Similar pages](#)

[Gatsby Course: Unsupervised Learning : Homepage](#)
Unsupervised Learning (Fall 2000). ... Syllabus (resources page): 10/10 1 - Introduction to **Unsupervised Learning** Geoff project: (ps, pdf). ...
www.gatsby.ucl.ac.uk/~quaid/course/ - 15k - [Cached](#) - [Similar pages](#)
[[More results from www.gatsby.ucl.ac.uk](#)]

[\[PDF\] Unsupervised Learning of the Morphology of a Natural Language](#)
File Format: PDF/Adobe Acrobat - [View as HTML](#)
Page 1. Page 2. Page 3. Page 4. Page 5. Page 6. Page 7. Page 8. Page 9. Page 10. Page 11. Page 12. Page 13. Page 14. Page 15. Page 16. Page 17. Page 18. Page 19 ...
acl.ldc.upenn.edu/J/J01/J01-2001.pdf - [Similar pages](#)

[Unsupervised Learning - The MIT Press](#)
... From Bradford Books: **Unsupervised Learning** Foundations of Neural Computation Edited by Geoffrey Hinton and Terrence J. Sejnowski Since its founding in 1989 by ...
mitpress.mit.edu/book-home.tcl?isbn=026258168X - 13k - [Cached](#) - [Similar pages](#)

[\[PS\] Unsupervised Learning of Disambiguation Rules for Part of](#)
File Format: Adobe PostScript - [View as Text](#)
Unsupervised Learning of Disambiguation Rules for Part of. Speech Tagging. Eric Brill. 1. ... It is possible to use **unsupervised learning** to train stochastic. ...
www.cs.jhu.edu/~brill/acl-wkshp.ps - [Similar pages](#)

[The Unsupervised Learning Group \(ULG\) at UT Austin](#)
The **Unsupervised Learning** Group (ULG). What ? The **Unsupervised Learning** Group (ULG) is a group of graduate students from the Computer ...
www.lans.ece.utexas.edu/ulg/ - 14k - [Cached](#) - [Similar pages](#)

 **Result Page:** [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [8](#) [9](#) [10](#) [Next](#)

Categorisation
Clustering
Relations between pages

Why a statistical approach?

- A probabilistic model of the data can be used to
 - make inferences about missing inputs
 - generate predictions/fantasies/imagery
 - make decisions which minimise expected loss
 - communicate the data in an efficient way
- Statistical modelling is equivalent to other views of learning:
 - information theoretic: finding compact representations of the data
 - physical analogies: minimising free energy of a corresponding statistical mechanical system

Basic Rules of Probability

Probabilities are non-negative $P(x) \geq 0 \forall x$.

Probabilities normalise: $\sum_{x \in \mathcal{X}} P(x) = 1$ for distributions if x is a discrete variable and $\int_{-\infty}^{+\infty} p(x) dx = 1$ for probability densities over continuous variables

The **joint probability** of x and y is: $P(x, y)$.

The **marginal probability** of x is: $P(x) = \sum_y P(x, y)$, assuming y is discrete.

The **conditional probability** of x given y is: $P(x|y) = P(x, y)/P(y)$

Bayes Rule:

$$P(x, y) = P(x)P(y|x) = P(y)P(x|y) \quad \Rightarrow$$

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

Warning: I will not be obsessively careful in my use of p and P for probability density and probability distribution. Should be obvious from context.

Information, Probability and Entropy

Information is the **reduction of uncertainty**. How do we measure uncertainty?

Some axioms (informal):

- if something is certain its uncertainty = 0
- uncertainty should be maximum if all choices are equally probable
- uncertainty (information) should add for independent sources

This leads to a discrete random variable X having uncertainty equal to the **entropy** function:

$$H(X) = - \sum_{X=x} P(X = x) \log P(X = x)$$

measured in *bits* (**binary digits**) if the base 2 logarithm is used or *nats* (**natural digits**) if the natural (base e) logarithm is used.

Some Definitions and Intuitions

- Surprise (for event $X = x$): $-\log P(X = x)$
- Entropy = average surprise: $H(X) = -\sum_{X=x} P(X = x) \log_2 P(X = x)$
- Conditional entropy

$$H(X|Y) = -\sum_x \sum_y P(x, y) \log_2 P(x|y)$$

- Mutual information

$$I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) = H(X) + H(Y) - H(X, Y)$$

- Kullback-Leibler divergence (relative entropy)

$$KL(P(X)||Q(X)) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

- Relation between Mutual information and KL: $I(X; Y) = KL(P(X, Y)||P(X)P(Y))$
- Independent random variables: $P(X, Y) = P(X)P(Y)$
- Conditional independence $X \perp\!\!\!\perp Y|Z$ (X conditionally independent of Y given Z)
means $P(X, Y|Z) = P(X|Z)P(Y|Z)$ and $P(X|Y, Z) = P(X|Z)$

Shannon's Source Coding Theorem

A discrete random variable X , distributed according to $P(X)$ has **entropy** equal to:

$$H(X) = - \sum_x P(x) \log P(x)$$

Shannon's source coding theorem: n independent samples of the random variable X , with entropy $H(X)$, can be compressed into minimum expected code of length $n\mathcal{L}$, where

$$H(X) \leq \mathcal{L} < H(X) + \frac{1}{n}$$

If each symbol is given a code length $l(x) = -\log_2 Q(x)$ then the expected per-symbol length \mathcal{L}_Q of the code is

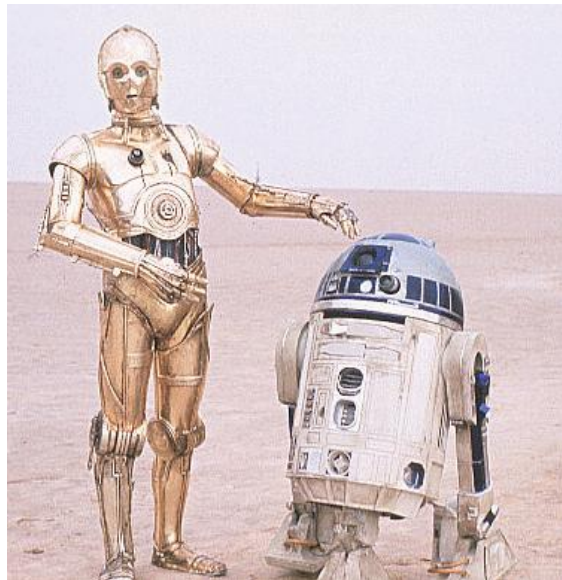
$$H(X) + KL(P\|Q) \leq \mathcal{L}_Q < H(X) + KL(P\|Q) + \frac{1}{n},$$

where the **relative-entropy** or **Kullback-Leibler divergence** is

$$KL(P\|Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)} \geq 0$$

Learning: A Statistical Approach II

- Goal: to represent the beliefs of learning agents.
- Cox Axioms lead to the following:
If plausibilities/beliefs are represented by real numbers, then the only reasonable and consistent way to manipulate them is Bayes rule.
- Frequency vs belief interpretation of probabilities
- The Dutch Book Theorem:
If you are willing to bet on your beliefs, then unless they satisfy Bayes rule there will always be a set of bets (“Dutch book”) that you would accept which is guaranteed to lose you money, no matter what outcome!



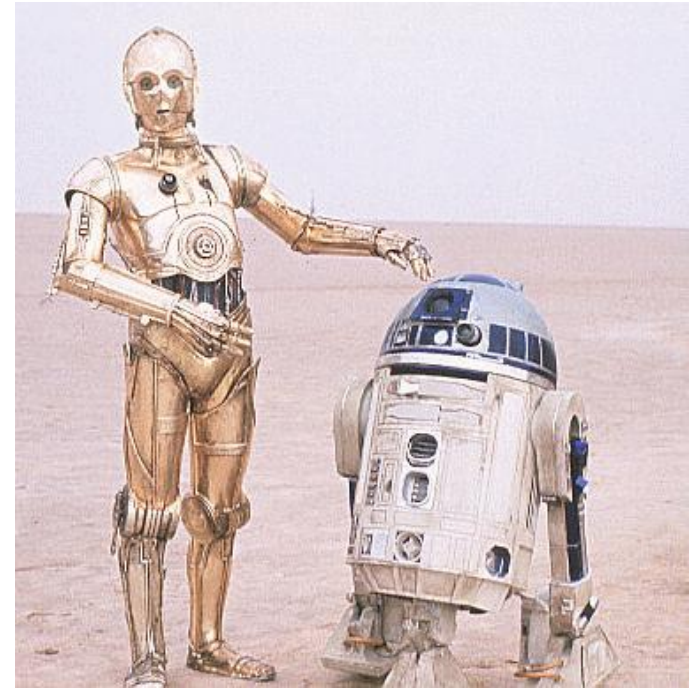
Representing Beliefs (Artificial Intelligence)

Consider a robot. In order to behave intelligently the robot should be able to represent beliefs about propositions in the world:

“my charging station is at location (x,y,z) ”

“my rangefinder is malfunctioning”

“that stormtrooper is hostile”



We want to represent the **strength** of these beliefs numerically in the brain of the robot, and we want to know what rules (calculus) we should use to manipulate those beliefs.

Representing Beliefs II

Let's use $b(x)$ to represent the strength of belief in (plausibility of) proposition x .

$$0 \leq b(x) \leq 1$$

$$b(x) = 0 \quad x \text{ is definitely **not true**}$$

$$b(x) = 1 \quad x \text{ is definitely **true**}$$

$$b(x|y) \quad \text{strength of belief that } x \text{ is true given that we know } y \text{ is true}$$

Cox Axioms (Desiderata):

- Strengths of belief (degrees of plausibility) are represented by real numbers
- Qualitative correspondence with common sense
- Consistency
 - If a conclusion can be reasoned in more than one way, then every way should lead to the same answer.
 - The robot always takes into account all relevant evidence.
 - Equivalent states of knowledge are represented by equivalent plausibility assignments.

Consequence: Belief functions (e.g. $b(x)$, $b(x|y)$, $b(x, y)$) must satisfy the rules of probability theory, including Bayes rule. (see Jaynes, *Probability Theory: The Logic of Science*)

The Dutch Book Theorem

Assume you are willing to accept bets with odds proportional to the strength of your beliefs. That is, $b(x) = 0.9$ implies that you will accept a bet:

$$x \text{ at } 1 : 9 \Rightarrow \begin{cases} x \text{ is true} & \text{win } \geq \text{£}1 \\ x \text{ is false} & \text{lose } \text{£}9 \end{cases}$$

Then, unless your beliefs satisfy the rules of probability theory, including Bayes rule, there exists a set of simultaneous bets (called a “Dutch Book”) which you are willing to accept, and for which **you are guaranteed to lose money, no matter what the outcome.**

E.g. suppose $A \cap B = \emptyset$, then

$$\left\{ \begin{array}{l} b(A) = 0.3 \\ b(B) = 0.2 \\ b(A \cup B) = 0.6 \end{array} \right\} \Rightarrow \text{accept the bets } \left\{ \begin{array}{l} \neg A \text{ at } 3 : 7 \\ \neg B \text{ at } 2 : 8 \\ A \cup B \text{ at } 4 : 6 \end{array} \right\}$$

But then:

$$\begin{aligned} \neg A \cap B &\Rightarrow \text{win } +3 - 8 + 4 = -1 \\ A \cap \neg B &\Rightarrow \text{win } -7 + 2 + 4 = -1 \\ \neg A \cap \neg B &\Rightarrow \text{win } +3 + 2 - 6 = -1 \end{aligned}$$

The only way to guard against Dutch Books is to ensure that your beliefs are coherent: i.e. satisfy the rules of probability.

Bayesian Learning

Apply the basic rules of probability to learning from data.

- Problem specification:

Data: $\mathcal{D} = \{x_1, \dots, x_n\}$ Models: $\mathcal{M}_1, \mathcal{M}_2$, etc. Parameters: θ_i (per model)

Prior probability of models: $P(\mathcal{M}_i)$.

Prior probabilities of model parameters: $P(\theta_i|\mathcal{M}_i)$

Model of data given parameters (likelihood model): $P(x|\theta_i, \mathcal{M}_i)$

- Data probability (**likelihood**)

$$P(\mathcal{D}|\theta_i, \mathcal{M}_i) = \prod_{j=1}^n P(x_j|\theta_i, \mathcal{M}_i) \equiv \mathcal{L}(\theta_i)$$

(provided the data are independently and identically distributed (**iid**).)

- Parameter learning (**posterior**):

$$P(\theta_i|\mathcal{D}, \mathcal{M}_i) = \frac{P(\mathcal{D}|\theta_i, \mathcal{M}_i)P(\theta_i|\mathcal{M}_i)}{P(\mathcal{D}|\mathcal{M}_i)}; \quad P(\mathcal{D}|\mathcal{M}_i) = \int d\theta_i P(\mathcal{D}|\theta_i, \mathcal{M}_i)P(\theta|\mathcal{M}_i)$$

$P(\mathcal{D}|\mathcal{M}_i)$ is called the **marginal likelihood** or **evidence** for \mathcal{M}_i . It is proportional to the posterior probability model \mathcal{M}_i being the one that generated the data.

- Model selection:

$$P(\mathcal{M}_i|\mathcal{D}) = \frac{P(\mathcal{D}|\mathcal{M}_i)P(\mathcal{M}_i)}{P(\mathcal{D})}$$

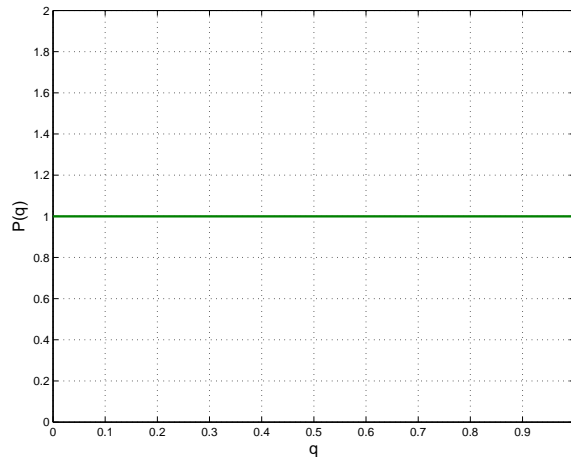
Bayesian Learning: A coin toss example

Coin toss: One parameter q — the odds of obtaining *heads*

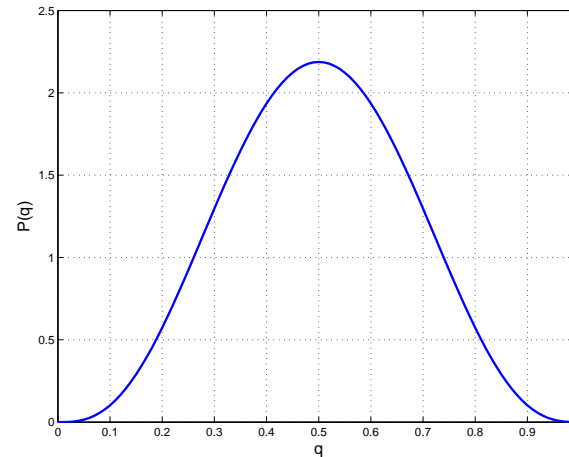
So our space of models is the set of distributions over $q \in [0, 1]$.

Learner A believes model \mathcal{M}_A : all values of q are equally plausible;

Learner B believes model \mathcal{M}_B : more plausible that the coin is “fair” ($q \approx 0.5$) than “biased”.



$$\mathbf{A}: \alpha_1 = \alpha_2 = 1.0$$



$$\mathbf{B}: \alpha_1 = \alpha_2 = 4.0$$

Both **prior beliefs** can be described by the Beta distribution:

$$p(q|\alpha_1, \alpha_2) = \frac{q^{(\alpha_1-1)}(1-q)^{(\alpha_2-1)}}{B(\alpha_1, \alpha_2)} = \text{Beta}(q|\alpha_1, \alpha_2)$$

Bayesian Learning: The coin toss (cont)

Now we observe a toss. Two possible outcomes:

$$p(\mathbf{H}|q) = q \quad p(\mathbf{T}|q) = 1 - q$$

Suppose our single coin toss comes out *heads*

The probability of the observed data (likelihood) is:

$$p(\mathbf{H}|q) = q$$

Using **Bayes Rule**, we multiply the prior, $p(q)$ by the likelihood and renormalise to get the posterior probability:

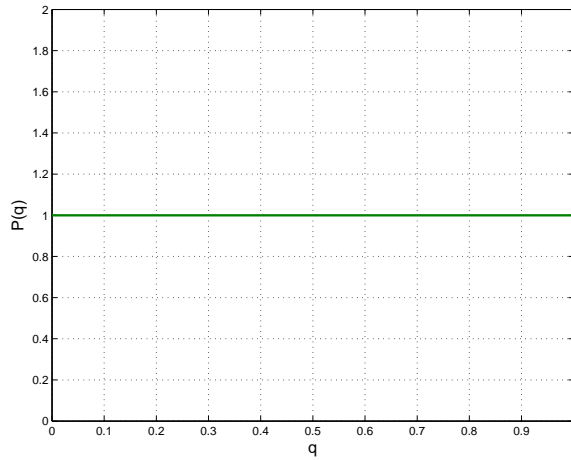
$$\begin{aligned} p(q|\mathbf{H}) &= \frac{p(q)p(\mathbf{H}|q)}{p(\mathbf{H})} \propto q \text{Beta}(q|\alpha_1, \alpha_2) \\ &\propto q q^{(\alpha_1-1)}(1-q)^{(\alpha_2-1)} = \text{Beta}(q|\alpha_1 + 1, \alpha_2) \end{aligned}$$

Bayesian Learning: The coin toss (cont)

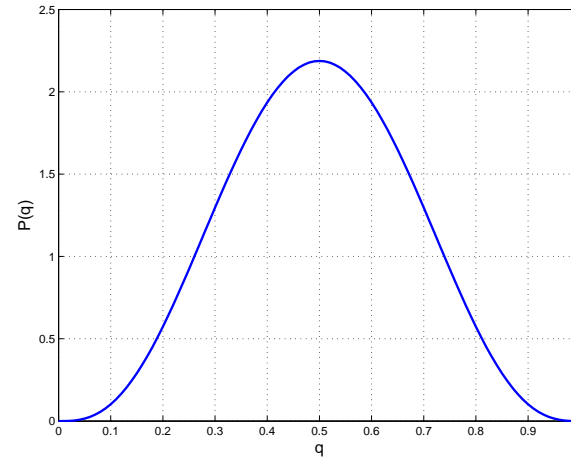
A

B

Prior

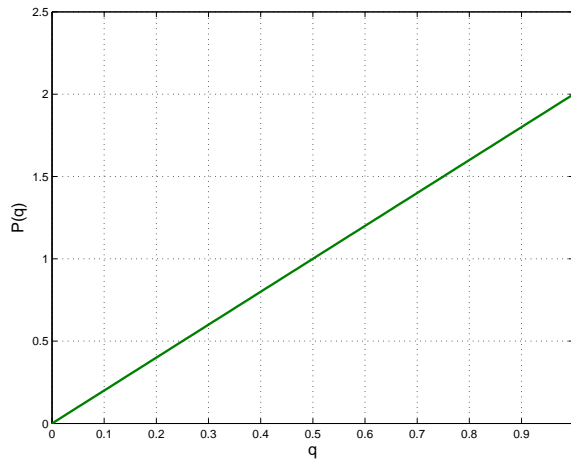


$\text{Beta}(q|1, 1)$

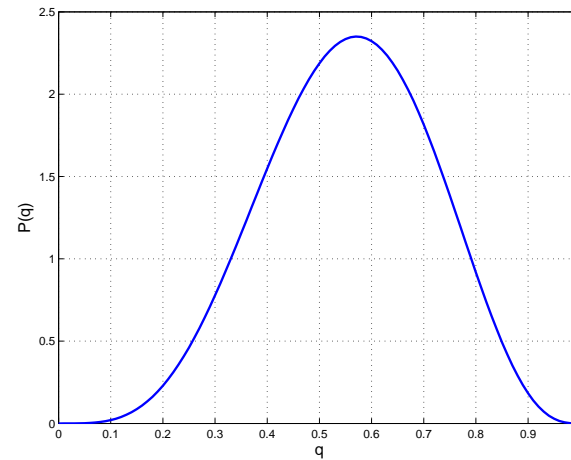


$\text{Beta}(q|4, 4)$

Posterior



$\text{Beta}(q|2, 1)$



$\text{Beta}(q|5, 4)$

Bayesian Learning: The coin toss (cont)

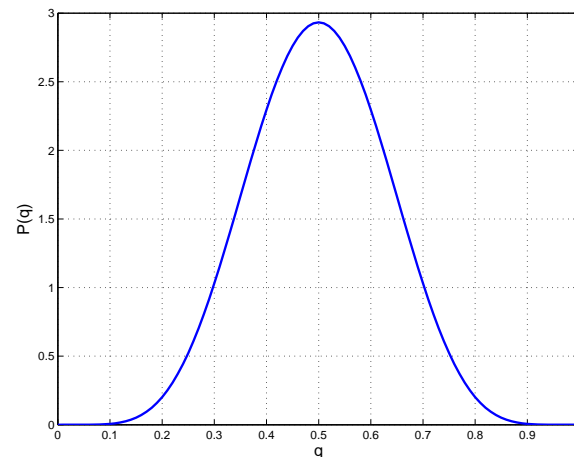
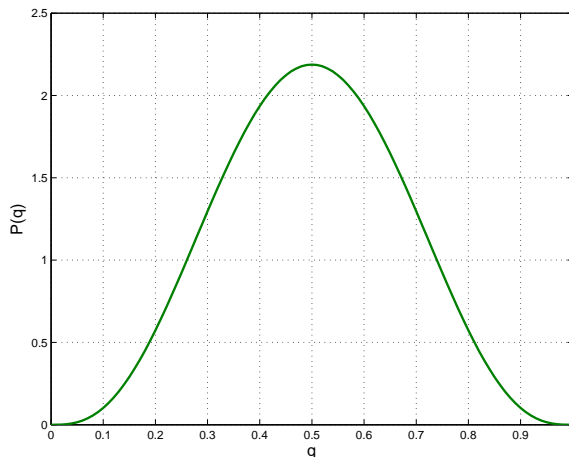
What about multiple tosses? Suppose we observe $\mathcal{D} = \{ H H T H T T \}$:

$$p(\{ H H T H T T \} | q) = qq(1 - q)q(1 - q)(1 - q) = q^3(1 - q)^3$$

This is still straightforward:

$$p(q | \mathcal{D}) = \frac{p(q)p(\mathcal{D}|q)}{p(\mathcal{D})} \propto q^3(1 - q)^3 \text{Beta}(q | \alpha_1, \alpha_2)$$

$$\propto \text{Beta}(q | \alpha_1 + 3, \alpha_2 + 3)$$



Conjugate priors

Updating the prior to form the posterior was particularly easy in these examples. This is because we used a **conjugate prior** for an **exponential family** likelihood.

Exponential family distributions take the form:

$$P(x|\theta) = g(\theta)f(x)e^{\phi(\theta)^\top \mathbf{T}(x)}$$

with $g(\theta)$ the normalising constant. Given n iid observations,

$$P(\{x_i\}|\theta) = \prod_i P(x_i|\theta) = g(\theta)^n e^{\phi(\theta)^\top \left(\sum_i \mathbf{T}(x_i)\right)} \prod_i f(x_i)$$

Thus, if the prior takes the **conjugate** form

$$P(\theta) = F(\boldsymbol{\tau}, \nu)g(\theta)^\nu e^{\phi(\theta)^\top \boldsymbol{\tau}}$$

with $F(\boldsymbol{\tau}, \nu)$ the normaliser, then the posterior is

$$P(\theta|\{x_i\}) \propto P(\{x_i\}|\theta)P(\theta) \propto g(\theta)^{\nu+n} e^{\phi(\theta)^\top \left(\boldsymbol{\tau} + \sum_i \mathbf{T}(x_i)\right)}$$

with the normaliser given by $F\left(\boldsymbol{\tau} + \sum_i \mathbf{T}(x_i), \nu + n\right)$.

Conjugate priors

The posterior given an exponential family likelihood and conjugate prior is:

$$P(\theta|\{x_i\}) = F(\boldsymbol{\tau} + \sum_i \mathbf{T}(x_i), \nu + n) g(\theta)^{\nu+n} \exp \left[\boldsymbol{\phi}(\theta)^\top \left(\boldsymbol{\tau} + \sum_i \mathbf{T}(x_i) \right) \right]$$

Here,

- $\boldsymbol{\phi}(\theta)$ is the vector of **natural parameters**
- $\sum_i \mathbf{T}(x_i)$ is the vector of **sufficient statistics**
- $\boldsymbol{\tau}$ are **pseudo-observations** which define the prior
- ν is the **scale** of the prior (need not be an integer)

As new data come in, each one increments the sufficient statistics vector and the scale to define the posterior.

The prior appears to be based on “pseudo-observations”, but:

1. This is different to applying Bayes’ rule. **No prior!** Sometimes we can take a uniform prior (say on $[0, 1]$ for q), but for unbounded θ , there may be no equivalent.
2. A valid conjugate prior might have non-integral ν or impossible $\boldsymbol{\tau}$, with no likelihood equivalent.

Conjugacy in the coin flip

Distributions are not always written in their natural exponential form.

The Bernoulli distribution (a single coin flip) with parameter q and observation $x \in \{0, 1\}$, can be written:

$$\begin{aligned}P(x|q) &= q^x(1 - q)^{(1-x)} \\&= e^{x \log q + (1-x) \log(1-q)} \\&= e^{\log(1-q) + x \log(q/(1-q))} \\&= (1 - q)e^{\log(q/(1-q))x}\end{aligned}$$

So the natural parameter is the **log odds** $\log(q/(1 - q))$, and the sufficient stats (for multiple tosses) is the number of heads.

The conjugate prior is

$$\begin{aligned}P(q) &= F(\tau, \nu) (1 - q)^\nu e^{\log(q/(1-q))\tau} \\&= F(\tau, \nu) (1 - q)^\nu e^{\tau \log q - \tau \log(1-q)} \\&= F(\tau, \nu) (1 - q)^{\nu - \tau} q^\tau\end{aligned}$$

which has the form of the Beta distribution $\Rightarrow F(\tau, \nu) = 1/B(\tau + 1, \nu - \tau + 1)$.

In general, then, the posterior will be $P(q|\{x_i\}) = \text{Beta}(\alpha_1, \alpha_2)$, with

$$\alpha_1 = 1 + \tau + \sum_i x_i \qquad \alpha_2 = 1 + (\nu + n) - \left(\tau + \sum_i x_i\right)$$

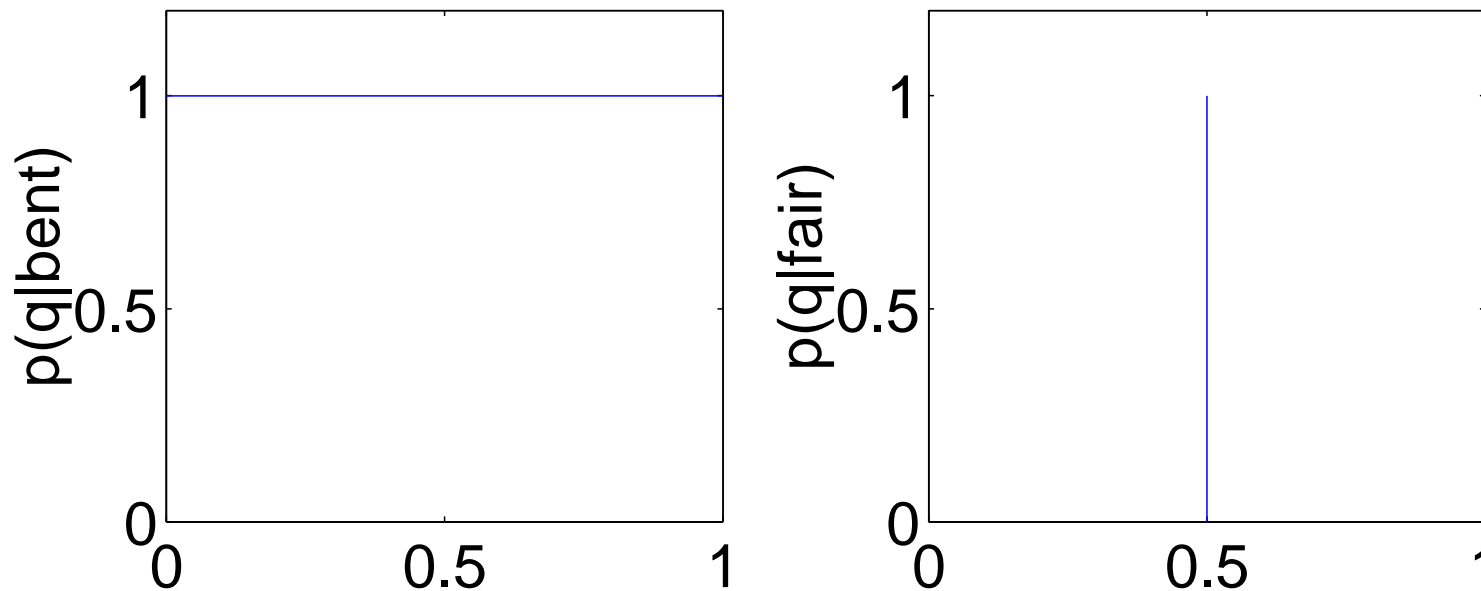
If we observe a head, we add 1 to the sufficient statistic $\sum x_i$, and also 1 to the count n . This increments α_1 . If we observe a tail we add 1 to n , but not to $\sum x_i$, incrementing α_2 .

Bayesian Coins II

We have seen how to update posteriors within each model. To study the choice of model, consider two more extreme models: “fair” and “bent”. A priori, we may think that “fair” is more probable, eg:

$$p(\text{fair}) = 0.8, \quad p(\text{bent}) = 0.2$$

For the bent coin, we assume all parameter values are equally likely, whilst the fair coin has a fixed probability:



We make 10 tosses, and get: $\mathcal{D} = (\text{T H T H T T T T T T})$.

Bayesian Coins II

Which model should we prefer *a posteriori* (i.e. after seeing the data)?

The **evidence** for the fair model is:

$$P(\mathcal{D}|\text{fair}) = (1/2)^{10} \approx 0.001$$

and for the bent model is:

$$P(\mathcal{D}|\text{bent}) = \int dq P(\mathcal{D}|q, \text{bent})p(q|\text{bent}) = \int dq q^2(1 - q)^8 = B(3, 9) \approx 0.002$$

Thus, the posterior for the models, by Bayes rule:

$$P(\text{fair}|\mathcal{D}) \propto 0.0008, \quad P(\text{bent}|\mathcal{D}) \propto 0.0004,$$

ie, a two-thirds probability that the coin is fair.

How do we make predictions? Could choose the fair model (model selection).

Or could weight the predictions from each model by their probability (model averaging).

Probability of H at next toss is:

$$P(H|\mathcal{D}) = P(H|\text{fair})P(\text{fair}|\mathcal{D}) + P(H|\text{bent})P(\text{bent}|\mathcal{D}) = \frac{2}{3} \times \frac{1}{2} + \frac{1}{3} \times \frac{3}{12} = \frac{5}{12}.$$

Parameter Learning

If an agent is learning parameters, it could report different aspects of the posterior (or likelihood).

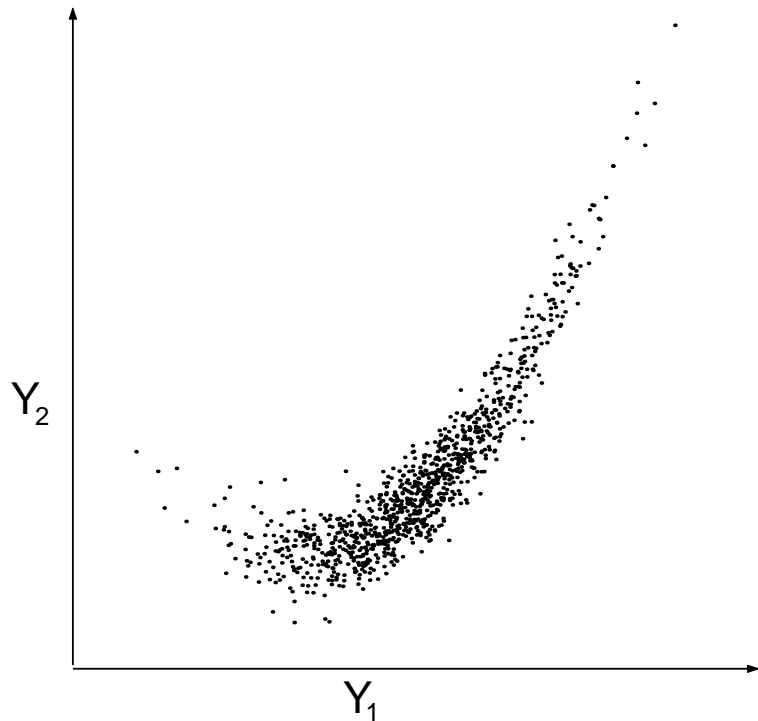
- **Bayesian Learning:** Assumes a prior over the model parameters. Computes the posterior distribution of the parameters: $P(\theta|\mathcal{D})$.
- **Maximum a Posteriori (MAP) Learning:** Assumes a prior over the model parameters $P(\theta)$. Finds a parameter setting that maximises the posterior: $P(\theta|\mathcal{D}) \propto P(\theta)P(\mathcal{D}|\theta)$.
- **Maximum Likelihood (ML) Learning:** Does not assume a prior over the model parameters. Finds a parameter setting that maximises the likelihood function: $P(\mathcal{D}|\theta)$.

Choosing between these and other alternatives may be a matter of definition, of goals, or of practicality

In practice (outside the exponential family), the Bayesian ideal may be computationally challenging, and may need to be approximated at best.

We will return to the Bayesian formulation in a while. For the next few weeks we will look at ML and MAP learning in more complex models.

Simple Statistical Modelling: modelling correlations



Assume:

- we have a data set $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- each data point is a vector of D features:
 $\mathbf{x}_i = [x_{i1} \dots x_{iD}]$
- the data points are i.i.d. (independent and identically distributed).

One of the simplest forms of unsupervised learning: model the **mean** of the data and the **correlations** between the D features in the data.

We can use a multivariate Gaussian model:

$$p(\mathbf{x}|\mu, \Sigma) = \mathcal{N}(\mu, \Sigma) = |2\pi\Sigma|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \mu)^\top \Sigma^{-1}(\mathbf{x} - \mu) \right\}$$

ML Estimation of a Gaussian

Data set $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, likelihood: $p(\mathcal{D}|\mu, \Sigma) = \prod_{n=1}^N p(\mathbf{x}_n|\mu, \Sigma)$

Goal: find μ and Σ that maximise likelihood \Leftrightarrow maximise log likelihood:

$$\begin{aligned}\ell &= \log \prod_{n=1}^N p(\mathbf{x}_n|\mu, \Sigma) = \sum_n \log p(\mathbf{x}_n|\mu, \Sigma) \\ &= -\frac{N}{2} \log |2\pi\Sigma| - \frac{1}{2} \sum_n (\mathbf{x}_n - \mu)^\top \Sigma^{-1} (\mathbf{x}_n - \mu)\end{aligned}$$

Note: equivalently, minimise $-\ell$, which is *quadratic* in μ

Procedure: take derivatives and set to zero:

$$\frac{\partial \ell}{\partial \mu} = 0 \quad \Rightarrow \quad \hat{\mu} = \frac{1}{N} \sum_n \mathbf{x}_n \quad (\text{sample mean})$$

$$\frac{\partial \ell}{\partial \Sigma} = 0 \quad \Rightarrow \quad \hat{\Sigma} = \frac{1}{N} \sum_n (\mathbf{x}_n - \hat{\mu})(\mathbf{x}_n - \hat{\mu})^\top \quad (\text{sample covariance})$$

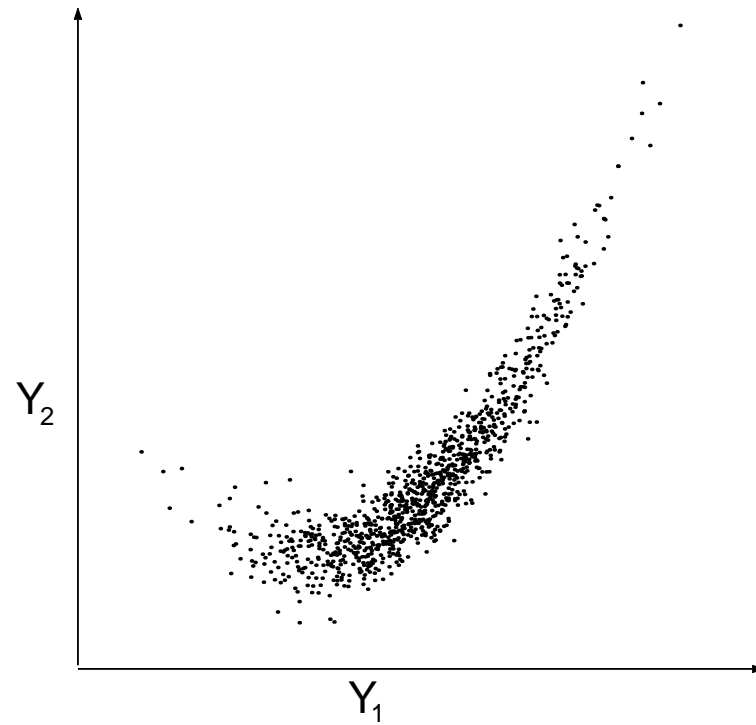
Gaussian Derivatives

$$\begin{aligned}\frac{\partial(-\ell)}{\partial\mu} &= \frac{\partial}{\partial\mu} \left[\frac{N}{2} \log |2\pi\Sigma| + \frac{1}{2} \sum_n (\mathbf{x}_n - \mu)^\top \Sigma^{-1} (\mathbf{x}_n - \mu) \right] \\ &= \frac{1}{2} \sum_n \frac{\partial}{\partial\mu} [(\mathbf{x}_n - \mu)^\top \Sigma^{-1} (\mathbf{x}_n - \mu)] \\ &= \frac{1}{2} \sum_n \frac{\partial}{\partial\mu} [\mathbf{x}_n^\top \Sigma^{-1} \mathbf{x}_n + \mu^\top \Sigma^{-1} \mu - 2\mu^\top \Sigma^{-1} \mathbf{x}_n] \\ &= \frac{1}{2} \sum_n \frac{\partial}{\partial\mu} [\mu^\top \Sigma^{-1} \mu] - 2 \frac{\partial}{\partial\mu} [\mu^\top \Sigma^{-1} \mathbf{x}_n] \\ &= \frac{1}{2} \sum_n [2\Sigma^{-1} \mu - 2\Sigma^{-1} \mathbf{x}_n] \\ &= N\Sigma^{-1} \mu - \Sigma^{-1} \sum_n \mathbf{x}_n \\ &= 0 \Rightarrow \hat{\mu} = \frac{1}{N} \sum_n \mathbf{x}_n\end{aligned}$$

Gaussian Derivatives

$$\begin{aligned}\frac{\partial(-\ell)}{\partial \Sigma^{-1}} &= \frac{\partial}{\partial \Sigma^{-1}} \left[\frac{N}{2} \log |2\pi \Sigma| + \frac{1}{2} \sum_n (\mathbf{x}_n - \mu)^\top \Sigma^{-1} (\mathbf{x}_n - \mu) \right] \\ &= \frac{\partial}{\partial \Sigma^{-1}} \left[\frac{N}{2} \log |2\pi I| \right] - \frac{\partial}{\partial \Sigma^{-1}} \left[\frac{N}{2} \log |\Sigma^{-1}| \right] \\ &\quad + \frac{1}{2} \sum_n \frac{\partial}{\partial \Sigma^{-1}} \left[(\mathbf{x}_n - \mu)^\top \Sigma^{-1} (\mathbf{x}_n - \mu) \right] \\ &= -\frac{N}{2} \Sigma^\top + \frac{1}{2} \sum_n (\mathbf{x}_n - \mu) (\mathbf{x}_n - \mu)^\top \\ = 0 &\Rightarrow \hat{\Sigma} = \frac{1}{N} \sum_n (\mathbf{x}_n - \mu) (\mathbf{x}_n - \mu)^\top\end{aligned}$$

Note



modelling correlations



maximising likelihood of a Gaussian model



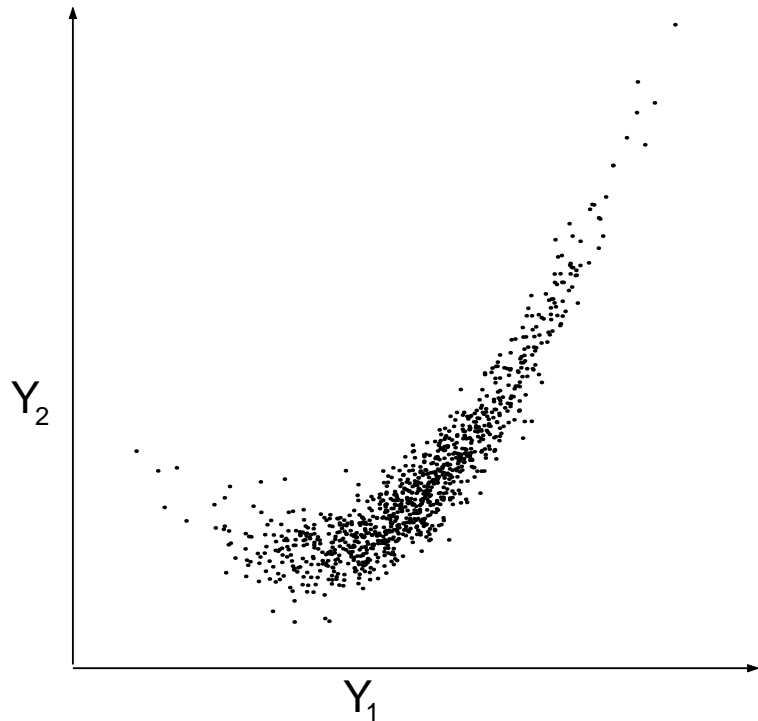
minimising a squared error cost function



minimizing data coding cost in bits (assuming Gaussian distributed)

Multivariate Linear Regression

The relationship between variables can also be modelled as a **conditional** distribution.



- data $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{y}_1) \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$
- each \mathbf{x}_i (\mathbf{y}_i) is a vector of D_x (D_y) features,
- \mathbf{y}_i is conditionally independent of all else, given \mathbf{x}_i .

One of the simplest forms of **supervised learning**: model \mathbf{y} as a **linear** function of \mathbf{x} , with **Gaussian** noise.

$$p(\mathbf{y}|\mathbf{x}, \mathbf{W}, \Sigma_y) = |2\pi\Sigma_y|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2}(\mathbf{y} - \mathbf{W}\mathbf{x})^\top \Sigma_y^{-1}(\mathbf{y} - \mathbf{W}\mathbf{x}) \right\}$$

Multivariate Linear Regression

ML estimates are obtained by maximising the (conditional) likelihood, as before:

$$\begin{aligned}\ell &= \sum_i \log p(\mathbf{y}_i | \mathbf{x}_i, \mathbf{W}, \Sigma_y) \\ &= -\frac{N}{2} \log |2\pi \Sigma_y| - \frac{1}{2} \sum_i (\mathbf{y}_i - \mathbf{W}\mathbf{x}_i)^\top \Sigma_y^{-1} (\mathbf{y}_i - \mathbf{W}\mathbf{x}_i)\end{aligned}$$

$$\begin{aligned}\frac{\partial(-\ell)}{\partial \mathbf{W}} &= \frac{\partial}{\partial \mathbf{W}} \left[\frac{N}{2} \log |2\pi \Sigma_y| + \frac{1}{2} \sum_i (\mathbf{y}_i - \mathbf{W}\mathbf{x}_i)^\top \Sigma_y^{-1} (\mathbf{y}_i - \mathbf{W}\mathbf{x}_i) \right] \\ &= \frac{1}{2} \sum_i \frac{\partial}{\partial \mathbf{W}} [(\mathbf{y}_i - \mathbf{W}\mathbf{x}_i)^\top \Sigma_y^{-1} (\mathbf{y}_i - \mathbf{W}\mathbf{x}_i)] \\ &= \frac{1}{2} \sum_i \frac{\partial}{\partial \mathbf{W}} [\mathbf{y}_i^\top \Sigma_y^{-1} \mathbf{y}_i + \mathbf{x}_i^\top \mathbf{W}^\top \Sigma_y^{-1} \mathbf{W}\mathbf{x}_i - 2\mathbf{x}_i^\top \mathbf{W}^\top \Sigma_y^{-1} \mathbf{y}_i] \\ &= \frac{1}{2} \sum_i \left[\frac{\partial}{\partial \mathbf{W}} \text{Tr} [\mathbf{W}^\top \Sigma_y^{-1} \mathbf{W}\mathbf{x}_i \mathbf{x}_i^\top] - 2 \frac{\partial}{\partial \mathbf{W}} \text{Tr} [\mathbf{W}^\top \Sigma_y^{-1} \mathbf{y}_i \mathbf{x}_i^\top] \right] \\ &= \frac{1}{2} \sum_i [2\Sigma_y^{-1} \mathbf{W}\mathbf{x}_i \mathbf{x}_i^\top - 2\Sigma_y^{-1} \mathbf{y}_i \mathbf{x}_i^\top] \\ &= 0 \Rightarrow \hat{\mathbf{W}} = \left(\sum_i \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \sum_i \mathbf{y}_i \mathbf{x}_i^\top\end{aligned}$$

Posterior estimation

Let y_i be scalar (so that W is a row vector) and write \mathbf{w} for the column vector of weights.

A conjugate prior for \mathbf{w} is

$$P(\mathbf{w}|\mathbf{A}) = \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1})$$

Then the **log** posterior on \mathbf{w} is

$$\begin{aligned}\log P(\mathbf{w}|\mathcal{D}, \mathbf{A}, \sigma_y) &= \log P(\mathcal{D}|\mathbf{w}, \mathbf{A}, \sigma_y) + \log P(\mathbf{w}|\mathbf{A}, \sigma_y) - \log P(\mathcal{D}|\mathbf{A}, \sigma_y) \\ &= -\frac{1}{2}\mathbf{w}^\top \mathbf{A} \mathbf{w} - \frac{1}{2} \sum_i (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 \sigma_y^{-2} + \text{const} \\ &= -\frac{1}{2} \mathbf{w}^\top \underbrace{\left(\mathbf{A} + \sigma_y^{-2} \sum_i \mathbf{x}_i \mathbf{x}_i^\top \right)}_{\Sigma_w^{-1}} \mathbf{w} + \mathbf{w}^\top \sum_i (y_i \mathbf{x}_i) \sigma_y^{-2} + \text{const} \\ &= -\frac{1}{2} \mathbf{w}^\top \Sigma_w^{-1} \mathbf{w} + \mathbf{w}^\top \underbrace{\Sigma_w^{-1} \Sigma_w}_{\mu_w} \sum_i (y_i \mathbf{x}_i) \sigma_y^{-2} + \text{const} \\ &= \log \mathcal{N}(\Sigma_w \sum_i (y_i \mathbf{x}_i) \sigma_y^{-2}, \Sigma_w)\end{aligned}$$

Gaussians for Regression

- Modelling the **conditional** $P(\mathbf{y}|\mathbf{x})$.
- If we also model $P(\mathbf{x})$, then learning is indistinguishable from **unsupervised**. In particular if $P(\mathbf{x})$ is Gaussian, and $P(\mathbf{y}|\mathbf{x})$ is linear-Gaussian, then \mathbf{x}, \mathbf{y} are **jointly Gaussian**.
- **Generalised Linear Models** (GLMs) generalise to non-Gaussian, exponential-family distributions and to non-linear **link functions**.

$$y_i \sim \text{ExpFam}(\mu_i, \phi)$$
$$g(\mu_i) = \mathbf{w}^\top \mathbf{x}_i$$

Posterior, or even ML, estimation is not possible in closed form \Rightarrow **iterative** methods such as gradient ascent or **iteratively re-weighted least squares (IRLS)**.

A warning to fMRIers: SPM uses GLM for “general” (not -ised) linear model; which is just linear.

- These models: Gaussians, Linear-Gaussian Regression and GLMs are important building blocks for the more sophisticated models we will develop later.
- Gaussian models are also used for regression in **Gaussian Process Models**. We’ll see these later too.

Three limitations of the multivariate Gaussian model

- What about higher order statistical structure in the data?

⇒ nonlinear and hierarchical models

- What happens if there are outliers?

⇒ other noise models

- There are $D(D + 1)/2$ parameters in the multivariate Gaussian model.
What if D is very large?

⇒ dimensionality reduction

End Notes

It is very important that you *understand* all the material in the following cribsheet:

<http://www.gatsby.ucl.ac.uk/teaching/courses/ul-2006/cribsheet.pdf>

The following notes by Sam Roweis are quite useful:

Matrix identities and matrix derivatives:

<http://www.cs.toronto.edu/~roweis/notes/matrixid.pdf>

Gaussian identities:

<http://www.cs.toronto.edu/~roweis/notes/gaussid.pdf>

Here is a useful statistics / pattern recognition glossary:

<http://research.microsoft.com/~minka/statlearn/glossary/>

Tom Minka's in-depth notes on matrix algebra:

<http://research.microsoft.com/~minka/papers/matrix/>