# COMPUTATIONAL STATISTICS
# PRIMER ON PROBABILITY AND STATISTICS

Luca Bortolussi

Department of Mathematics and Geosciences
University of Trieste

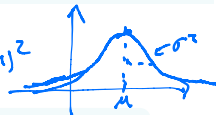Office 303, third floor, E1 3
luca@dmi.units.it

Trieste, Winter Semester 2015/2016

# Outline

# MULTIVARIATE NORMAL DISTRIBUTION

*[handwritten annotations: 1D; $p(x|\mu,\sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$]*

*[handwritten: $=E[x_i]$; $\mu_i = \langle x_i \rangle$; $\Sigma_{ij} = cov(x_i, x_j)$]*

This is the most important distribution we will use, and generalises the 1d normal. In $d$ dimensions

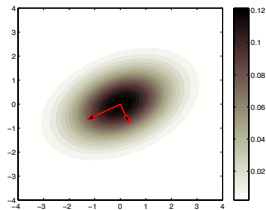*[handwritten: d-dim vector; $\Sigma$ symm + pos.defined; d×d; quadratic form]*

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

It holds $\boldsymbol{\mu} = \langle \mathbf{x} \rangle$, and $\boldsymbol{\Sigma} = \mathrm{cov}(\mathbf{x}, \mathbf{x}) = \langle (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \rangle$

(a)

(b)

# PROPERTIES OF MULTIVARIATE NORMAL

**Completing the square**

A useful technique in manipulating Gaussians is completing the square. For example, the expression

$$\exp\left(-\frac{1}{2}\mathbf{x}^\mathsf{T}\mathbf{A}\mathbf{x} + \mathbf{b}^\mathsf{T}\mathbf{x}\right) \qquad (8.4.10)$$

can be transformed as follows. First we complete the square:

$$\frac{1}{2}\mathbf{x}^\mathsf{T}\mathbf{A}\mathbf{x} - \mathbf{b}^\mathsf{T}\mathbf{x} = \frac{1}{2}\left(\mathbf{x} - \mathbf{A}^{-1}\mathbf{b}\right)^\mathsf{T}\mathbf{A}\left(\mathbf{x} - \mathbf{A}^{-1}\mathbf{b}\right) - \frac{1}{2}\mathbf{b}^\mathsf{T}\mathbf{A}^{-1}\mathbf{b} \qquad (8.4.11)$$

Hence

$$\exp\left(-\frac{1}{2}\mathbf{x}^\mathsf{T}\mathbf{A}\mathbf{x} + \mathbf{b}^\mathsf{T}\mathbf{x}\right) = \mathcal{N}\left(\mathbf{x}|\mathbf{A}^{-1}\mathbf{b}, \mathbf{A}^{-1}\right)\sqrt{\det\left(2\pi\mathbf{A}^{-1}\right)}\exp\left(\frac{1}{2}\mathbf{b}^\mathsf{T}\mathbf{A}^{-1}\mathbf{b}\right) \qquad (8.4.12)$$

$p(\mathbf{x}|\mathbf{A},\mathbf{b}) c \exp\left(-\frac{1}{2}\mathbf{x}\mathbf{A}\mathbf{x}^\mathsf{T} + \mathbf{b}^\mathsf{T}\mathbf{x}\right)$ is known as the canonical representation, and it is normal with mean $\mathbf{A}^{-1}\mathbf{b}$ and covariance $\mathbf{A}^{-1}$.
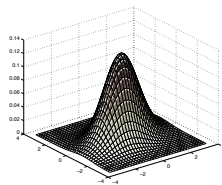
**Linear transformation**

**Result 8.3** (Linear Transform of a Gaussian). Let $\mathbf{y}$ be linearly related to $\mathbf{x}$ through

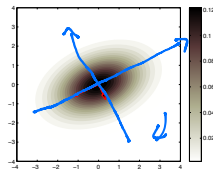$$\mathbf{y} = \mathbf{M}\mathbf{x} + \boldsymbol{\eta} \qquad (8.4.14)$$

where $\mathbf{x} \perp\!\!\!\perp \boldsymbol{\eta}$, $\boldsymbol{\eta} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, and $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x)$. Then the marginal $p(\mathbf{y}) = \int_\mathbf{x} p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$ is a Gaussian

$$p(\mathbf{y}) = \mathcal{N}\left(\mathbf{y}|\mathbf{M}\boldsymbol{\mu}_x + \boldsymbol{\mu}, \mathbf{M}\boldsymbol{\Sigma}_x\mathbf{M}^\mathsf{T} + \boldsymbol{\Sigma}\right) \qquad (8.4.15)$$

# PROPERTIES OF MULTIVARIATE NORMAL



(a)                              (b)

Eigendecomposition

$$\boldsymbol{\Sigma} = \mathbf{E}\boldsymbol{\Lambda}\mathbf{E}^{\mathsf{T}} \tag{8.4.5}$$

where $\mathbf{E}^{\mathsf{T}}\mathbf{E} = \mathbf{I}$ and $\boldsymbol{\Lambda} = \operatorname{diag}(\lambda_1, \ldots, \lambda_D)$. In the case of a covariance matrix, all the eigenvalues $\lambda_i$ are positive. This means that one can use the transformation

$$\mathbf{y} = \boldsymbol{\Lambda}^{-\frac{1}{2}}\mathbf{E}^{\mathsf{T}}(\mathbf{x} - \boldsymbol{\mu}) \tag{8.4.6}$$

so that

$$(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}}\boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}}\mathbf{E}\boldsymbol{\Lambda}^{-1}\mathbf{E}^{\mathsf{T}}(\mathbf{x} - \boldsymbol{\mu}) = \mathbf{y}^{\mathsf{T}}\mathbf{y} \tag{8.4.7}$$

So by rescaling, we can obtain a product of *d*-univariate standard normal distributions, one per dimension.

# PROPERTIES OF MULTIVARIATE NORMAL

Marginal and conditional of multivariate Gaussians

**Result 8.4** (Partitioned Gaussian). Consider a distribution $\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ defined jointly over two vectors $\mathbf{x}$ and $\mathbf{y}$ of potentially differing dimensions,

$$\mathbf{z} = \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \qquad (8.4.16)$$

with corresponding mean and partitioned covariance

$$\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{pmatrix} \qquad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{pmatrix} \qquad (8.4.17)$$

where $\boldsymbol{\Sigma}_{yx} \equiv \boldsymbol{\Sigma}_{xy}^{\mathsf{T}}$. The marginal distribution is given by

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_{xx}) \qquad (8.4.18)$$

$$p(x) = \int p(x,y)\,dy$$

and conditional

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}\left(\mathbf{x}|\boldsymbol{\mu}_x + \boldsymbol{\Sigma}_{xy}\boldsymbol{\Sigma}_{yy}^{-1}(\mathbf{y} - \boldsymbol{\mu}_y), \boldsymbol{\Sigma}_{xx} - \boldsymbol{\Sigma}_{xy}\boldsymbol{\Sigma}_{yy}^{-1}\boldsymbol{\Sigma}_{yx}\right) \qquad (8.4.19)$$

# PROPERTIES OF MULTIVARIATE NORMAL

Product of multivariate Gaussians

**Result 8.2** (Product of two Gaussians). The product of two Gaussians is another Gaussian, with a multiplicative factor, exercise(8.35):

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)\, \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})\, \frac{\exp\left(-\frac{1}{2}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^{\mathsf{T}} \mathbf{S}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)\right)}{\sqrt{\det(2\pi \mathbf{S})}} \tag{8.4.8}$$

where $\mathbf{S} \equiv \boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2$ and the mean and covariance are given by

$$\boldsymbol{\mu} = \boldsymbol{\Sigma}_1 \mathbf{S}^{-1} \boldsymbol{\mu}_2 + \boldsymbol{\Sigma}_2 \mathbf{S}^{-1} \boldsymbol{\mu}_1 \qquad \boldsymbol{\Sigma} = \boldsymbol{\Sigma}_1 \mathbf{S}^{-1} \boldsymbol{\Sigma}_2 \tag{8.4.9}$$
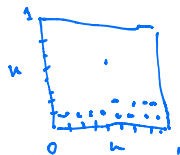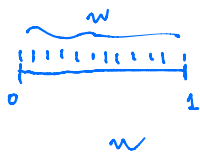
Gaussian average of a quadratic function

**Result 8.5** (Gaussian average of a quadratic function).

$$\left\langle \mathbf{x}^{\mathsf{T}} \mathbf{A} \mathbf{x} \right\rangle_{\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})} = \boldsymbol{\mu}^{\mathsf{T}} \mathbf{A} \boldsymbol{\mu} + \operatorname{trace}(\mathbf{A}\boldsymbol{\Sigma}) \quad \Leftarrow$$

# THE CURSE OF DIMENSIONALITY

*Exercise I*: Suppose you want to explore uniformly a region by gridding it. How many grid points do you need?

# THE CURSE OF DIMENSIONALITY

*Exercise II*: Suppose you sample from a uniform distribution in $d$ dimensions. What is the probability of finding a point inside the region $[\epsilon, 1 - \epsilon]^d$?

$$U = (U_1, \ , U_d) \qquad U_i \sim UNIFORM(0,1)$$

$$\vec{u} \in [\epsilon, 1-\epsilon]^d \iff u_i \in [\epsilon, 1-\epsilon] \qquad (P(u_i^2)=1)$$

$$\underbrace{\phantom{[\epsilon, 1-\epsilon]^d}}_{R_\epsilon} \qquad P(u_i \in [\epsilon, 1-\epsilon]) = 1-2\epsilon$$

$$\forall \epsilon, \quad P(\vec{u} \in R_\epsilon) = (1-2\epsilon)^d \xrightarrow[d \to \infty]{} 0$$

# THE CURSE OF DIMENSIONALITY

*d-dim, indep. standard gauss.*

*Exercise III*: Suppose you sample from a spherical Gaussian distribution. Where do the points lie as the dimensions increase?

# MIXTURES: HOW TO BUILD MORE DISTRIBUTIONS

- More general distributions can be built via mixtures: e.g.

$$p(x|\mu_{1...,n}, \sigma^2_{1,...,n}) = \sum_i \pi_i \mathcal{N}(\mu_i, \sigma^2_i)$$

$$[0,1] \quad \sum \pi_i = 1$$

where the *mixing coefficients* $\pi_i$ are discretely distributed

# MIXTURES: HOW TO BUILD MORE DISTRIBUTIONS

- More general distributions can be built via mixtures: e.g.

$$p(x|\mu_{1...,n}, \sigma^2_{1,...,n}) = \sum_i \pi_i \mathcal{N}(\mu_i, \sigma^2_i)$$

where the *mixing coefficients* $\pi_i$ are discretely distributed

- You can interpret this as a two stage hierarchical process: choose one component out of a discrete distribution, then choose the distribution for that component

# MIXTURES: HOW TO BUILD MORE DISTRIBUTIONS

- **IMPORTANT CONCEPT**: the mixture

$$p(x|\mu_{1...,n}, \sigma^2_{1,...,n}) = \sum_i \pi_i \mathcal{N}(\mu_i, \sigma^2_i)$$

is an example of *latent variable model*, with a latent class variable and an observed continuous value. The mixture is the marginal distribution for the observations (w.r.t. the latent variable)

$$P(x) = \sum_\pi P(x, \pi) = \sum_\pi P(\pi) \cdot P(x|\pi)$$

$$p(\pi) = \sum_x p(\pi, x) = \sum_x p(\pi|x) p(x)$$

$\vec{x}$ oss. $p(\pi | \vec{x}) \sim$ Bayes.

# MIXTURES: HOW TO BUILD MORE DISTRIBUTIONS

- **IMPORTANT CONCEPT**: the mixture

$$p(x|\mu_{1...,n}, \sigma^2_{1,...,n}) = \sum_i \pi_i \mathcal{N}(\mu_i, \sigma^2_i)$$

  is an example of *latent variable model*, with a latent class variable and an observed continuous value. The mixture is the marginal distribution for the observations (w.r.t. the latent variable)
- The probability of the latent variables given the observations can be obtained using Bayes' theorem.

# CONTINUOUS MIXTURES: SOME COOL DISTRIBUTIONS

- No need for the mixing distribution (latent variable) to be discrete

# CONTINUOUS MIXTURES: SOME COOL DISTRIBUTIONS

- No need for the mixing distribution (latent variable) to be discrete
- Suppose you are interested in the means of normally distributed samples (possibly with different variances/ precisions): Marginalising the precision in a Gaussian using a Gamma mixing distribution yields a *Student t-distribution*

# CONTINUOUS MIXTURES: SOME COOL DISTRIBUTIONS

- No need for the mixing distribution (latent variable) to be discrete
- Suppose you are interested in the means of normally distributed samples (possibly with different variances/ precisions): Marginalising the precision in a Gaussian using a Gamma mixing distribution yields a *Student t-distribution*
- Suppose you have multiple rare event processes happening with slightly different rates: Marginalising the rate in a Poisson distribution using a Gamma mixing distribution yields a *negative binomial* distribution

# Outline

# PARAMETERS?

- Many distributions are written as conditional probabilities *given* the parameters: $p(x|\theta)$
- Often the values of the parameters are not known
- If we have observations, we can try to estimate the parameters from such data.
- We assume to have independent and identically distributed (i.i.d.) observations of $p(x|\theta_{true})$: $\mathbf{x} = x_1, \ldots, x_N$.

# MAXIMUM LIKELIHOOD

- Likelihood for i.i.d. observations $\mathbf{x} = x_1, \ldots, x_N$:

$$p(\mathbf{x}|\theta) = \prod_{i=1}^{N} p(x_i|\theta)$$

independence

# MAXIMUM LIKELIHOOD

- Likelihood for i.i.d. observations $\mathbf{x} = x_1, \ldots, x_N$:

$$p(\mathbf{x}|\theta) = \prod_{i=1}^{N} p(x_i|\theta)$$

- Choose the parameters that best explain the observations: we pick $\theta$ by maximum likelihood:

$$\hat{\theta} = \operatorname{argmax}_{\theta} \left[ \prod_i p(x_i|\theta) \right]$$

## MAXIMUM A POSTERIORI

- Suppose we can encode prior knowledge (or absence of it) in a prior distribution over parameters, $p(\theta)$.

- We can then compute the posterior distribution, given i.i.d. observations $\mathbf{x} = x_1, \ldots, x_N$, by Bayes theorem:

$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})}$$

EVIDENCE
MARGINAL LIKELIHOOD

where

$$p(\mathbf{x}) = \int_{\theta} p(\mathbf{x}|\theta)p(\theta)d\theta$$

## MAXIMUM A POSTERIORI

- Suppose we can encode prior knowledge (or absence of it) in a prior distribution over parameters, $p(\theta)$.
- We can then compute the posterior distribution, given i.i.d. observations $\mathbf{x} = x_1, \ldots, x_N$, by Bayes theorem:

$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})}$$

where

$$p(\mathbf{x}) = \int_\theta p(\mathbf{x}|\theta)p(\theta)d\theta$$

- Estimate $\theta_{true}$ by the *maximum a posteriori (MAP)* estimate

$$\hat{\theta}_{MAP} = \mathrm{argmax}_\theta \left[ p(\theta) \prod_i p(x_i|\theta) \right]$$

# EXERCISE: FITTING A DISCRETE DISTRIBUTION

- We have a discrete distribution with values in $K = \{1, \ldots, k\}$, with parameters $\boldsymbol{\mu} = \mu_1, \ldots, \mu_k$, $\sum_i \mu_i = 1$.
- We have independent observations $\mathrm{x} = x_1, \ldots, x_N$, each taking values in $K$.
- The likelihood is

$$\mathcal{L}(\boldsymbol{\mu}) = p(\mathrm{x}|\boldsymbol{\mu}) = \prod_{i=1}^{N} p(x_i|\boldsymbol{\mu})$$

- Compute the Maximum Likelihood estimate of $\boldsymbol{\mu}$. What is the intuitive meaning of the result? What happens if one of the $k$ values is not represented in your sample?

## EXERCISE: FITTING A DISCRETE DISTRIBUTION

$$x_i = (0 \cdots 1 \cdots 0) \sim (x_{i1}, \cdots x_{ik}) \qquad \sum_j x_{ij} = 1$$

$$p(x_i | \vec{\mu}) = \prod_j \mu_j^{x_{ij}} \qquad \left( \text{se } x_{i\bar{k}} \in S, \text{ allora } p(x_i | \vec{\mu}) = \mu_S \right)$$

$$p(\vec{x} | \vec{\mu}) = \prod_{i,j} \mu_j^{x_{ij}} = \prod_j \mu_j^{n_j} \qquad n_j = \# (j, \vec{x})$$

$$\ell(\vec{\mu}) = \log p(\vec{x} | \vec{\mu}) = \sum_j n_j \log \mu_j$$

$$\begin{cases} \text{Impone } \sum_j \mu_j = 1 \\ \lambda - \text{molt. di forap.} \end{cases}$$

$$\ell(\vec{\mu}, \hat{\lambda}) = \sum_j n_j \log \mu_j + \lambda \left( \sum_j \mu_j - 1 \right)$$

$$\frac{\partial \ell}{\partial \mu_i} = \frac{n_j}{\mu_j} + \lambda = 0 \qquad \frac{\partial \ell}{\partial \lambda} = \sum_j \mu_j - 1 = 0$$

$$\mu_j = -\frac{n_j}{\lambda}$$

$$\lambda = -N \qquad \boxed{\hat{\mu}_j = \frac{n_j}{N}}$$

# EXERCISE II: FITTING A GAUSSIAN DISTRIBUTION

We have independent, real valued observations $x = x_1, \ldots, x_N$.
Fit a Gaussian by maximum likelihood.

$$\prod_i p(x_i \mid \mu, \sigma^2) = \prod_i \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x_i - \mu)^2}$$

$$\ell(\mu, \sigma^2) = -\frac{N}{2}\log 2\pi\sigma^2 - \frac{1}{2\sigma^2}\sum_i (x_i - \mu)^2$$

$$\frac{\partial \ell}{\partial \mu} = 0 \qquad \sum_i (x_i - \mu) = 0 \qquad \hat{\mu} = \frac{\sum_i x_i}{N}$$

$$\frac{\partial \ell}{\partial \sigma^2} = 0 \qquad -\frac{N}{2} \cdot \frac{1}{\sigma^2} + \frac{1}{2}\sum_i (x_i - \mu)^2 \cdot \frac{1}{(\sigma^2)^2} = 0 \qquad \hat{\sigma}^2 = \frac{\sum_i (x_i - \hat{\mu})^2}{N}$$

# BAYESIAN ESTIMATION

- The Bayesian approach fully quantifies uncertainty
- The parameters are treated as additional random variables with their own *prior* distribution $p(\theta)$
- The observation likelihood is combined with the prior to obtain a *posterior* distribution via Bayes' theorem

$$p(\theta|x) = \frac{p(x|\theta)p(\theta)}{p(x)}$$

$$P(x) = \int P(x|\theta)P(\theta)d\theta$$

- The distribution of the observable $x$ (*predictive* distribution) is obtained as

$$p(x|x) = \int p(x|\theta)p(\theta|x)d\theta$$

# EXERCISE: BAYESIAN FITTING OF GAUSSIANS

- Let data $x_i \quad i = 1, \dots, N$ be distributed according to a Gaussian with mean $\mu$ and variance $\sigma^2$
- Let the prior distribution over the mean $\mu$ be a Gaussian with mean $m$ and variance $v^2$
- Compute the posterior (and predictive distribution, Exercise) $\uparrow$ GAUSSIANA $\mathcal{N}(x \mid \mu_P, \sigma_P^2)$

$$\mu_P = m \cdot \frac{\sigma^2}{N v^2 + \sigma^2} + \hat{\mu} \cdot \frac{N v^2}{N v^2 + \sigma^2}$$

$$\sigma_P^2 = \left( \frac{1}{v^2} + \frac{N}{\sigma^2} \right)^{-1}$$

# Exercise: Bayesian fitting of Gaussians

## ESTIMATORS

- A procedure to calculate an expectation is called an *estimator*
- e.g., fitting a Gaussian to data by maximum likelihood provides the M.L. estimator for mean and variance, or Bayesian posterior mean
- An estimator will be a noisy estimate of the true value, due to finite sample effects
- An estimator $\hat{t}$ is *unbiased* if its expectation (under the joint distribution of the data set) coincides with the true value
- An estimator $\hat{t}$ is *consistent* if it converges to the true value when the number of data goes to infinity.

# EXERCISE: BIASED ESTIMATOR

The ML estimator of variance, $\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^{N} (x_i - \hat{\mu})^2$ is biased:

$\langle \hat{\sigma}^2 \rangle = \frac{N-1}{N} \sigma^2$.

$\uparrow \quad P(\vec{x})$

$\left[ \sigma^2 = \langle x_i^2 \rangle - \mu^2 \right.$

$\left\{ \begin{array}{l} \langle x_i \rangle = \mu \\ \langle x_i^2 \rangle = \sigma^2 + \mu^2 \\ \langle x_i x_j \rangle = \mu^2 \end{array} \right\}$

## BOOTSTRAP

- For an estimator, in theory we can compute its mean and its variance under the joint distribution of the datasets. In practice, getting the variance may be very hard. Bootstrapping can be used instead.
- Given the dataset $\mathrm{x} = x_1, \ldots, x_N$, construct from it $K$ new datasets $\mathrm{x}_i$, also of size $N$, by sampling with repetitions.
- compute the estimator $\hat{\theta}_i$ for each $\mathrm{x}_i$.
- Compute the empirical variance (or any other statistics) from $\mathrm{x}_1, \ldots, \mathrm{x}_K$.
- This is an estimate of the actual variance of the estimator.

# CONJUGATE PRIORS

- The Bayesian way has advantages in that it quantifies uncertainty and regularizes naturally
- BUT computing the normalisation in Bayes theorem is very hard
- The case when it is possible is when the prior and the posterior are of the same form (*conjugate*)
- Example + Exercise: Bernoulli and Beta.
- Example: discrete and Dirichlet
- Exercise: conjugate priors for the univariate normal (mean)

## CONJUGATE PRIORS: BERNOULLI BINOMIAL AND BETA

Show that the Beta is the conjugate prior for the Bernoulli distribution.

$$\text{BERNOULLI} \quad x \in \{0,1\} \quad P(x\mid\vartheta) = \vartheta^x(1-\vartheta)^{1-x}$$

$$x_1,-,x_N \quad P(\vec{x}\mid\vartheta) = \prod_1 P(x_i\mid\vartheta) = \vartheta^k(1-\vartheta)^{N-k}$$

$$k = \sum_i x_i = \#(1,\vec{x})$$

$$\text{Beta}(\vartheta\mid a,b) = C\cdot\vartheta^{a-1}(1-\vartheta)^{b-1}$$

$$P(\vartheta\mid\vec{x}) = \frac{1}{P(\vec{x})}\cdot\vartheta^k(1-\vartheta)^{N-k}\cdot C\,\vartheta^{a-1}(1-\vartheta)^{b-1}$$

$$= C'\,\vartheta^{a+k-1}(1-\vartheta)^{b+N-k-1} = \text{Beta}(\vartheta\mid a+k-1, b+N-k-1)$$

# Outline

# ENTROPY

- Probability theory is the basis of information theory (interesting, but not the topic of this course).
- An important quantity is the *entropy* of a distribution

# ENTROPY

- Probability theory is the basis of information theory (interesting, but not the topic of this course).
- An important quantity is the *entropy* of a distribution

$$H[p] = -\sum_i p_i \log_2 p_i$$

  Or for continuous distributions:

$$H[p] = -\int p(x) \log p(x) dx$$

- Entropy measures the level of disorder of a distribution; for discrete distributions, it is always $\geq 0$ and 0 only for deterministic distributions. The maximum is $\log K$, if $K$ is the size of the support of the discrete distribution, and is achieved by the uniform distribution.

## DIVERGENCE

- The *relative entropy* or *Kullback-Leibler (KL) divergence* between two distributions is

$$KL[q\|p] = \sum_i q_i \log \frac{q_i}{p_i} = \langle \log q \rangle_q - \langle \log p \rangle_q$$
$$H[q]$$

Of in the continuous case

$$KL[q\|p] = \int q(x) \log \frac{q(x)}{p(x)} dx$$

- Fact: *KL* is convex and $\geq 0$ (by Jensen ineq)
- Fact: *KL* is zero if and only if $p = q$.

## CONDITIONAL ENTROPY AND MUTUAL INFORMATION

- Conditional entropy is defined as

$$H[p(x|y)] = -\int\int p(x,y)\log p(x|y)dxdy = H[p(x,y)]-H[p(y)]$$

and captures the residual uncertainty on $x$ once $y$ is known.

- Mutual information between r.v. $x$ and $y$ is defined as

$$I[x,y] = KL[p(x,y)|p(x)p(y)] = H[p(x)] - H[p(x|y)]$$

and captures the reduction in uncertainty about $x$ by knowing $y$, i.e. it is a measure of how much $y$ brings information about $x$, and viceversa.

# JUSTIFICATION FOR MAXIMUM LIKELIHOOD

$x_i \in \{0, 1\}$

- Given a data set $\mathrm{x} = \{x_i\}, \quad i = 1, \dots, N$, let the empirical distribution be

$$p_{emp}(x) = \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(x_i)$$

$P_{emp}(0) = \frac{n_0}{N}$

$P_{emp}(1) = \frac{n_1}{N}$

with $\mathbb{I}$ the indicator function of a set

- To find a suitable distribution $q$ to model the data, one may wish to minimize the Kullback-Leibler divergence

$log$ -likelih.

$$KL[p_{emp} \| q] = H[p_{emp}] - \langle \log q(x) \rangle_{p_{emp}} = -\frac{1}{N} \sum \log q(x_i) \uparrow cost$$

- *Maximum likelihood is equivalent to minimizing a KL divergence with the empirical distirbution*

# Outline

ORARIO

LUN 16-18 (T)
MAR 14-15.30 (4)
MER 14-17.00 (L)

1 Basics of probability theory

2 Some probability distributions

3 Fitting distributions

( LUN 19 10 -> ore
                  16.30 )

4 Information theory

5 Decision theory

## AN OVERVIEW

- Suppose we have a classification problem, and we are able to learn a model of the joint distribution $p(x, y)$, where $y$ is a categorical variable. Given a new input $x^*$, for which we want to make a prediction, to which class should we assign it?

## AN OVERVIEW

- Suppose we have a classification problem, and we are able to learn a model of the joint distribution $p(x, y)$, where $y$ is a categorical variable. Given a new input $x^*$, for which we want to make a prediction, to which class should we assign it?

- We may choose to assign it to class $j$ if $p(y = j|x^*)$ is the maximum one. However, suppose $y$ models having or not a cancer, and that
$p(y = 0|x^*) = 0.51 > 0.49 = p(y = 1|x^*).$

## AN OVERVIEW

- Suppose we have a classification problem, and we are able to learn a model of the joint distribution $p(x, y)$, where $y$ is a categorical variable. Given a new input $x^*$, for which we want to make a prediction, to which class should we assign it?

- We may choose to assign it to class $j$ if $p(y = j|x^*)$ is the maximum one. However, suppose $y$ models having or not a cancer, and that
  $p(y = 0|x^*) = 0.51 > 0.49 = p(y = 1|x^*)$.

- To be more flexible, we can specify a loss function (or utility function), which is the cost $c_{k,j}$ of assigning $x^*$ to class $j$ when the true class is $k$.

## AN OVERVIEW

- Suppose we have a classification problem, and we are able to learn a model of the joint distribution $p(x, y)$, where $y$ is a categorical variable. Given a new input $x^*$, for which we want to make a prediction, to which class should we assign it?

- We may choose to assign it to class $j$ if $p(y = j|x^*)$ is the maximum one. However, suppose $y$ models having or not a cancer, and that
$p(y = 0|x^*) = 0.51 > 0.49 = p(y = 1|x^*)$.

- To be more flexible, we can specify a loss function (or utility function), which is the cost $c_{k,j}$ of assigning $x^*$ to class $j$ when the true class is $k$.

- Then we can assign a point $x^*$ to the class $j$ minimising the expected loss w.r.t. the learned joint distribution (i.e. $\sum_k c_{k,j} p(y = k|x^*)$).