Note: These scribed notes have only been lightly proofread.

## 9.1 Naive Bayes

### 9.1.1 Introduction

**Remarque** : Contrary to its name, "Naive Bayes" is not a Bayesian method

Let's Consider the following problem of classification $x \in \mathbb{X}^p \longmapsto y \in \{1, 2, \ldots, M\}$.

Here, $x = (x_1, x_2, \ldots, x_p)$ is a vector of descriptors (or features) : $\forall i \in \{1, 2, \ldots, p\}, x_i \in \mathbb{X}$, with $\mathbb{X} = \{1, 2, \ldots, K\}$ (or $\mathbb{X} = \mathbb{R}$).

Goal : Learning $p(y|x)$

A very naive method will trigger off a combinatorial explosion : $\theta \in \mathbb{R}^{K^p}$.

Bayes formula gets us :

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

The Naive Bayes method consists in assuming that the features $x_i$ are all conditionally independent from the class, hence :

$$p(x|y) = \prod_{i=1}^{p} p(x_i|y)$$

Then, the Bayes formula gives us:

$$p(y|x) = \frac{p(y) \prod_{i=1}^{p} p(x_i|y)}{p(x)} = \frac{p(y) \prod_{i=1}^{p} p(x_i|y)}{\sum_{y'} p(y') \prod_{i=1}^{p} p(x_i|y')}$$

We consider the case where the features take discrete values. Consequently the new graphical model contains only discrete random variables. Then, we can write a discrete model as an exponential family. Indeed we can write:

$$\log p(x_i = k|y = k') = \delta(x_i = k,\, y = k')\,\theta_{ikk'}$$

and

$$\log p(y = k') = \delta(y = k')\,\theta_{k'}$$

We can see that the dummy functions $\delta(x_i = k, y = k')$ and $\delta(y = k')$ are the *sufficient statistics* of the joint distribution model for $y$ and the variables $x_i$, where $\theta_{ikk'}$ and $\theta_{k'}$ are *canonical parameters*. Thus , we can write:

$$\log p(y, x_1, \dots, x_p) = \sum_{i,k,k'} \delta(x_i = k, y = k')\theta_{ikk'} + \sum_{k'} \delta(y = k')\theta_{k'} - A((\theta_{ikk'})_{i,k,k'}, (\theta_{k'})_{k'})$$

Where $A((\theta_{ikk'})_{i,k,k'}, (\theta_{k'})_{k'})$ is the log-partition function.

We have rewritten the joint distribution model of $(y, x_1, \dots, x_p)$ as an exponential family. Given that the maximum of likelihood estimator of an exponential family, where the canonical parameters are not combined, is also the maximum entropy estimator; as seen in a previous course and provided that the statistical moments of the sufficient statistics equal their empirical moments.

Thus, if we introduce

$$N_{ikk'} = \# \left\{ (x_i, y) = (k, k') \right\}$$
$$N = \sum_{i,k,k'} N_{ikk'},$$

The maximum likelihood estimator must satisfy the moment constraints

$$\widehat{p}(y = k') = \frac{\sum_{i,k} N_{ikk'}}{N} \qquad \text{et} \qquad \widehat{p}(x_i = k | y = k') = \frac{N_{ikk'}}{\sum_{k''} N_{ik''k'}},$$

which define them completly.

Then, we can write the estimators of the canonical parameters as:

$$\widehat{\theta}_{ikk'} = \log \widehat{p}(x_i = k | y = k') \qquad \text{et} \qquad \widehat{\theta}_{k'} = \log \widehat{p}(y = k').$$

However, our goal is to obtain a classification model, that is to say, a model of only the conditional probability law. From the approximated generative model and applying the Bayes rule we can get:

$$\log \widehat{p}(y = k' | x) = \sum_{i=1}^{p} \log \widehat{p}(x_i | y = k') + \log \widehat{p}(y = k') - \log \sum_{k'} \left( \widehat{p}(y = k') \prod_{i=1}^{p} \widehat{p}(x_i | y = k') \right)$$

We can re write the conditional model as an exponential family

$$\log p(y | x) = \sum_{i,k,k'} \delta(x_i = k, y = k')\theta_{ikk'} + \sum_{k'} \delta(y = k')\theta_{k'} - \log p(x)$$

Its sufficient statistics and canonical parameters are equal to those of the generative model, but seen as functions of the random variable $y$, given that $x$ is fixed (we could write $\phi_{x,i,k,k'}(y) = \delta(x_i = k, y = k')$). As for the log-partition function, it is now equal to $\log p(x)$.

Warning: $\widehat{\theta}_{ikk'}$ is the maximum likelihood estimator in the generative model which, usually, is not equal to the maximum likelihood estimator in the conditional model.

### 9.1.2   Advantages and Drawbacks

Advantages :

- Doable in line.

- Computationally tractable solution.

Drawbacks :

- Generative : generative models produce good estimator whenever the model is "true", or in statistical words *well specified*, which means that the process that generate the real data induce a distribution equal to the one of the generative model. When the model is not *well specified* (which is the most common case) we'd better use a discriminative method.

### 9.1.3   Discriminative method

The problem that we have considered in the previous section is the generative model for classification in K classes. How to learn, in a discriminatory way , a classifier in K classes? Is it possible to use an exponential family?

We have already seen the logistic regression for 2 classes classification:

$$p\left(y=1|x\right)=\frac{\exp\left(\omega^{T}x\right)}{1+\exp\left(\omega^{T}x\right)}$$

Let's study the K-multiclass logistic regression:

$$
\begin{aligned}
p\left(y=k'|x\right)&=\frac{\exp\left(\sum_{i=1}^{p}\sum_{k=1}^{K}\delta\left(x_i=k\right)\theta_{ikk'}\right)}{\sum_{k''=1}^{M}\exp\left(\sum_{i=1}^{p}\sum_{k=1}^{K}\delta\left(x_i=k\right)\theta_{ikk''}\right)}\\
&=\exp\left(\sum_{i=1}^{p}\sum_{k=1}^{K}\delta\left(x_i=k\right)\theta_{ikk'}-\log\left(\sum_{k''=1}^{M}\exp\left(\sum_{i=1}^{p}\sum_{k=1}^{K}\delta\left(x_i=k\right)\theta_{ikk''}\right)\right)\right)\\
&=\exp\left(\theta_{k'}^{T}\phi\left(x\right)-\log\left(\sum_{k''=1}^{M}\exp\left(\theta_{k''}^{T}\phi\left(x\right)\right)\right)\right)\\
&=\frac{\exp\left(\theta_{k'}^{T}\phi\left(x\right)\right)}{\sum_{k''=1}^{M}\exp\left(\theta_{k''}^{T}\phi\left(x\right)\right)}
\end{aligned}
$$

Although we have built the model from different staring consideration, the resulting modelling ( that is the set of possible distribution) is of the same exponential family than the Naive Bayes model.

Nonetheless, the fitted model in a discriminatory approach will be different from the one fitted in a generative approach : the fitting of the K-multiclass logistic regression results from the maximisation of the likelihood of the classes $y^{(j)}$ of a set of learning, given that $x^{(j)}$ are fixed. In other words, the fitting is obtained by computing the maximum likelihood estimator in the conditional model. Unlike what happens in the generative model, the estimator can't be obtained in a analytical form and the learning requires solving a numerical optimisation problem.

## 9.2   Bayesian Method

### 9.2.1   Introduction

Vocabulary :

- a priori : $p\,(\theta)$

- likelihood : $p\,(x|\theta)$

- marginal likelihood : $\int p\,(x|\theta)\,p\,(\theta)\,d\theta$

- a posteriori : $p\,(\theta|x)$

The Bayesian formulation enables us to introduce the a priori information in the process of estimation. For instance , let's imagine that we play heads or tails :

- with an "unknown" coin, we've got the information a priori : we'll use the uniform law for $p\,(\theta)$.

- with a "normal" coin , we'll use a distribution with an important concentration of mass around 0,5 for $p\,(\theta)$.

For a Bayesian, offering a "limited" estimator, as the maximum likelihood estimator, which gives a unique value for $\theta$, is not enough because the estimator itself do not translate the inherent uncertainty of the learning process. Thus, its estimator will be the density a posteriori, obtained from the Bayes rule, which is written in continuous notations as :

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

The Bayesian specifies the uncertainty with distributions that form its estimator, rather than combining an estimator with confidence intervals.

If the Bayesian is forced to produce a limited estimator, he uses the expectation of the underlying quantity under the a posteriori distribution; for instance for $\theta$:

$$\mu_{post} = \mathbb{E}\,[\theta|D] = \mathbb{E}\,[\theta|x_1, x_2, \ldots, x_n] = \int \theta p\,(\theta|x_1, x_2, \ldots, x_n)\,d\theta$$

### 9.2.2   a posteriori Maximum (PAM)

$$\theta_{MAP} = \arg\max_{\theta} p\,(\theta|x_1, x_2, \ldots, x_n)$$
$$= \arg\max_{\theta} p\,(x_1, x_2, \ldots, x_n|\theta)\,p\,(\theta)$$

Because, with the Bayes rule:

$$p\left(\theta|x_1, x_2, \ldots, x_n\right) = \frac{p\left(x_1, x_2, \ldots, x_n|\theta\right)p\left(\theta\right)}{p\left(x\right)}$$

The a posteriori maximum is not really Bayesian, it's rather a slight modification brought to the frequentist estimator.

### 9.2.3 Predictive probability

In the Bayesian paradigm, the probability of a future observation $x^*$ will be estimated by the *Predictive probability*:

$$
\begin{aligned}
p\left(x^*|D\right) &= p\left(x^*|x_1, x_2, \ldots, x_n\right) \\
&= \int p\left(x^*|\theta\right)p\left(\theta|x_1, x_2, \ldots, x_n\right)d\theta
\end{aligned}
$$

$$
\begin{aligned}
p\left(\theta|x_1, x_2, \ldots, x_n\right) &\propto p\left(x_n|\theta\right)p\left(x_1|\theta\right)p\left(x_2|\theta\right)\ldots p\left(x_{n-1}|\theta\right)p\left(\theta\right) \\
&\propto p\left(x_n|\theta\right)p\left(\theta|x_1, x_2, \ldots, x_{n-1}\right)p\left(x_1, x_2, \ldots, x_{n-1}\right) \\
&\propto p\left(x_n|\theta\right)p\left(\theta|x_1, x_2, \ldots, x_{n-1}\right)\frac{p\left(x_1, x_2, \ldots, x_{n-1}\right)}{p\left(x_1, x_2, \ldots, x_n\right)}
\end{aligned}
$$

A sequential calculus is possible since:

$$p\left(\theta|x_1, x_2, \ldots, x_n\right) = \frac{p\left(x_n|\theta\right)p\left(\theta|x_1, x_2, \ldots, x_{n-1}\right)}{p\left(x_n|x_1, x_2, \ldots, x_{n-1}\right)}$$

Vocabulary :

- a priori information : $p\left(\theta|x_1, x_2, \ldots, x_{n-1}\right)$

- likelihood : $p\left(x_n|\theta\right)$

- a posteriori information : $p\left(\theta|x_1, x_2, \ldots, x_n\right)$

$$p\left(x_1, x_2, \ldots, x_n\right) = \int \prod_{i=1}^{n} p\left(x_i|\theta\right)p\left(\theta\right)d\theta$$

### 9.2.4 Exchangeable situations

#### 9.2.4.1 Exchangeablility

The random variables $X_1$, $X_2$, ..., $X_n$ are exchangeable if they have the same distribution as $X_{\widehat{\Sigma}(1)}, X_{\widehat{\Sigma}(2)}, \ldots, X_{\widehat{\Sigma}(n)}$ for any permutation of indices $\widehat{\Sigma}$.

#### 9.2.4.2    de Finetti's theorem

Si $X_1$, $X_2$, ..., $X_n$ are exchangeable, then it exists a stochastic process $G$ such that :

$$p\left(x_1, x_2, \ldots, x_n\right) = \int \prod_{i=1}^{n} p\left(x_i|G\right) d\mu(G)$$

Where $d\mu(G)$ is the generalisation of "$p\left(G\right) dG$" for a stochastic process.

#### 9.2.4.3    Why do we care about exchangeable situations?

The i.i.d variables are a particular case of the situation of exchangeable variables, that we see in practice. However when the i.i.d data are combined with non scalar observations, the different components are no longer independent. In some cases, those components are nonetheless exchangeable. For instance in a text, words are shown as sequences that are not exchangeable because of the syntax. But if we forget the order of the words as in the "bag of word" model, then the components are exchangeable. It's the basic principle used in the LDA model.

### 9.2.5    Example of model

#### 9.2.5.1    Bernoulli variable

let's consider random variables $X_i \in \{0, 1\}$. We'll assume that the $X_i$ i.i.d. conditionally to $\theta$.

$$p\left(x|\theta\right) = \theta^x \left(1 - \theta\right)^{1-x}$$

#### 9.2.5.2    Priors

let's introduce the *distribution* Beta whose density on $[0, 1]$ is

$$p(\theta; \alpha, \beta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

Where $B(\alpha, \beta)$ is the alias of the Beta *function* :

$$\forall \alpha > 0, \forall \beta > 0, B\left(\alpha, \beta\right) = \int_0^1 \theta^{\alpha-1} \left(1 - \theta\right)^{\beta-1} d\theta$$

And the Gamma function :

$$\Gamma\left(x\right) = \int_0^{+\infty} t^{x-1} \exp\left(-t\right) dt$$

We can show that :

$$B\left(\alpha,\beta\right)=\frac{\Gamma\left(\alpha\right)\Gamma\left(\beta\right)}{\Gamma\left(\alpha+\beta\right)}$$

We choose as the prior distribution on $\theta$ the Beta distribution:

$$p\left(\theta\right)\propto\theta^{\alpha-1}\left(1-\theta\right)^{\beta-1}$$

$$p\left(\theta\right)=\frac{\theta^{\alpha-1}\left(1-\theta\right)^{\beta-1}}{B\left(\alpha,\beta\right)}$$

### 9.2.5.3    A posteriori

$$p\left(\theta|x\right)=\frac{p\left(x,\theta\right)}{p\left(x\right)}\propto p\left(x,\theta\right)$$

But :

$$p\left(x,\theta\right)=\theta^{x}\left(1-\theta\right)^{1-x}\frac{\theta^{\alpha-1}\left(1-\theta\right)^{\beta-1}}{B\left(\alpha,\beta\right)}$$

Hence :

$$p\left(\theta|x\right)\propto\frac{\theta^{x+\alpha-1}\left(1-\theta\right)^{1-x+\beta-1}}{B\left(\alpha,\beta\right)}$$

$$p\left(\theta|x\right)=\frac{\theta^{x+\alpha-1}\left(1-\theta\right)^{1-x+\beta-1}}{B\left(x+\alpha,1-x+\beta\right)}$$

Thus, if instead of considering a unique variable , we observe an i.i.d sample of data, the joint distribution can be written as :

$$\theta^{\alpha-1}\left(1-\theta\right)^{\beta-1}\prod_{i=1}^{n}\theta^{x_i}\left(1-\theta\right)^{1-x_i}.$$

Let's introduce :

$$k=\sum_{i=1}^{n}x_i$$

Then we get :

$$p\left(\theta|x_1,x_2,\ldots,x_n\right)=\frac{\theta^{k+\alpha-1}\left(1-\theta\right)^{n-k+\beta-1}}{B\left(k+\alpha,n-k+\beta\right)}$$

### 9.2.6   Distributions

$$\theta \sim Beta\left(\alpha, \beta\right)$$

For $\alpha = \beta = 1$, we've got a uniform prior.
For $\alpha = \beta > 1$, we've got a bell curve.
For $\alpha = \beta < 1$, we've got a U curve.
$\mathbb{E}\left[\theta\right] = \frac{\alpha}{\alpha+\beta}$
$\mathbb{V}\left[\theta\right] = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)} = \frac{\alpha}{(\alpha+\beta)} \times \frac{\beta}{(\alpha+\beta)} \times \frac{1}{(\alpha+\beta+1)}$
For $\alpha > 1$ and $\beta > 1$, we've got the mode : $\frac{\alpha-1}{\alpha+\beta-2}$.
In the case, let's write $D$ for the data:

$$\theta_{post} = \mathbb{E}\left[\theta|D\right] = \frac{\alpha + k}{\alpha + \beta + n} = \frac{\alpha}{(\alpha+\beta)} \times \frac{(\alpha+\beta)}{(\alpha+\beta+n)} + \frac{n}{(\alpha+\beta+n)} \times \frac{k}{n}$$

We can see that the a posteriori expectation of the parameter is a convex combination of the maximum likelihood estimator and the prior expectation. It converges asymptotically to the maximum likelihood estimator .

If we use a uniform prior distribution, $\mathbb{E}|\mathbb{D}\left[\theta\right] = \frac{k+1}{n+2}$. Laplace proposed to correct the frequentist estimator, it seemed odd to him that he was not defined in the absence of data. He proposed to add two virtual observation (0 and 1) such that in the absence of data the estimator equals $\frac{1}{2}$. This correction is known as *Laplace's correction*.

The variance of the a posteriori distribution decrease in $\frac{1}{n}$ .

$$\mathbb{V}\left[\theta|D\right] = \theta_M\left(1 - \theta_M\right)\frac{1}{(\alpha + \beta + n)}$$

We have chosen a sharper distribution around $\theta_M$, in the same way than in a frequentist approach, the confidence intervals narrow around the estimator when the number of observations increase.

### 9.2.7   Playful propriety

$$p\left(x_1, x_2, \ldots, x_n\right) = \frac{B\left(k + \alpha, n - k + \beta\right)}{B\left(\alpha, \beta\right)} = \frac{\Gamma\left(\alpha + k\right)\Gamma\left(\beta + n - k\right)\Gamma\left(\alpha + \beta\right)}{\Gamma\left(\alpha + \beta + n\right)\Gamma\left(\alpha\right)\Gamma\left(\beta\right)} \qquad (9.1)$$

let's use the propriety of the Gamma function:

$$\Gamma\left(n + 1\right) = n!$$

$$\text{and} \qquad \forall x > -1, \, \Gamma\left(x + 1\right) = x\Gamma\left(x\right)$$

such that

$$\Gamma\left(\alpha + k\right) = \left(\alpha + k - 1\right)\left(\alpha + k - 2\right)\ldots\alpha\Gamma\left(\alpha\right)$$

let's write $\alpha^{[k]} = \alpha\left(\alpha + 1\right)\ldots\left(\alpha + k - 1\right)$ and simplify the expression 9.1 :

$$p(x_1, x_2, \ldots, x_n) = \frac{\alpha^{[k]} \beta^{[n-k]}}{(\alpha + \beta)^{[n]}}$$

We shall note the analogy with the Polya urn model: let us consider $(\alpha + \beta)$ balls of colour : $\alpha$ are black, $\beta$ are white. When drawing a first black ball, the probability of the event is :

$$\mathbb{P}(X_1 = 1) = \frac{\alpha}{\alpha + \beta}$$

After the drawing, we put back the ball in the urn and we add a ball of the same colour. Let's imagine that we draw again a black ball then the probability of this event is:

$$\mathbb{P}(X_1 = 1, X_2 = 1) = \mathbb{P}(X_1 = 1)\mathbb{P}(X_2 = 1|X_1 = 1) = \frac{\alpha}{\alpha + \beta} \times \frac{\alpha + 1}{\alpha + \beta + 1}$$

However :

$$\mathbb{P}(X_1 = 1, X_2 = 0) = \frac{\alpha}{\alpha + \beta} \times \frac{\beta}{\alpha + \beta + 1}$$

In more general case , we show by recurrence that the marginal probability of obtaining some sequence of colours by drawing from a Polya urn is exactly the marginal probability of obtaining the same result from the marginal model, obtained by integrating on a priori *theta*. First, this show that drawings from a Polya urn are exchangeable; Secondly, the mechanism of this type of urn, and its exchangeability, we'll be useful for the Gibbs sampling and for the same type of Bayesian models.

## 9.2.8 Conjugate priors

Let $\mathbb{F}$ be a set. We assume that $p(x|\theta)$ known, we deduce from that : $p(\theta) \in \mathbb{F}$ such taht $p(\theta|x) \in \mathbb{F}$. We say that $p(\theta)$ is conjugated to the model $p(x|\theta)$.

### 9.2.8.1 Exponential model

Let's consider:

$$p(x|\theta) = \exp\left(\langle \theta, \phi(x) \rangle - A(\theta)\right)$$
$$p(\theta) = \exp\left(\langle \alpha, \theta \rangle - \tau A(\theta) - B(\alpha, \tau)\right)$$

For $p(x|\theta)$, $\theta$ is the canonical parameter. For $p(\theta)$, $\alpha$ is the canonical parameter and $\theta$ is the sufficient statistic. Let us note that $B$ do not stand for the Beta distribution.

$$p(\theta|x) \propto p(x|\theta)p(\theta) \propto \exp\left(\langle \theta, \phi(x) \rangle - A(\theta) + \langle \alpha, \theta \rangle - \tau A(\theta) - B(\alpha, \tau)\right)$$

Let us define :

$$\bar{\phi} = \frac{1}{n}\sum_{i=1}^{n} \phi(x_i)$$

Then :

$$p\left(\theta|x_i\right) \propto \exp\left(\langle\theta, \alpha + \phi\left(x_i\right)\rangle - \left(\tau + 1\right)A\left(\theta\right) - B\left(\alpha + \phi\left(x_i\right), \tau + 1\right)\right)$$

$$p\left(\theta|x_1, x_2, \ldots, x_n\right) \propto \exp\left(\langle\theta, \alpha + n\bar{\phi}\rangle - \left(\tau + n\right)A\left(\theta\right) - B\left(\alpha + n\bar{\phi}, \tau + n\right)\right)$$

$$p\left(x_1, x_2, \ldots, x_n\right) \propto \exp\left(B\left(\alpha, \tau\right) - B\left(\alpha + n\bar{\phi}, \tau + n\right)\right)$$

Since the family is an exponential one,

$$\nu_{post} = \mathbb{E}\left[\theta|D\right] = \nabla_\alpha B\left(\alpha + n\bar{\phi}, \tau + n\right)$$

$\theta_{MAP}$ results from :

$$\nabla_\theta p\left(\theta|x_1, x_2, \ldots, x_n\right) = 0$$
$$\alpha + n\bar{\phi} = \left(\tau + n\right)\nabla_\theta A\left(\theta\right) = \left(\tau + n\right)\mu\left(\theta\right)$$

Thus we get $\mu_{MAP} = \mu\left(\theta\right)$ in the previous equation. Consequently:

$$\mu_{MAP} = \frac{\alpha + n\bar{\phi}}{\tau + n} = \frac{\alpha}{\tau} \times \frac{\tau}{\tau + n} + \frac{n}{\tau + n}\bar{\phi}$$

### 9.2.8.2    Univariate Gaussian

**With and a priori on $\mu$ but not on $\sigma^2$**

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

$$p\left(\mu|\mu_0, \tau^2\right) = \frac{1}{\sqrt{2\pi\tau^2}}\exp\left(-\frac{1}{2}\frac{\left(\mu - \mu_0\right)^2}{\tau^2}\right)$$

Thus :

$$p\left(D|\mu, \sigma^2\right) = p\left(x_1, x_2, \ldots, x_n|\mu, \sigma^2\right)$$
$$= \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(-\frac{1}{2}\sum_{i=1}^{n}\frac{\left(x_i - \mu\right)^2}{\sigma^2}\right)$$

$$
\begin{aligned}
p\left(\mu|D\right) &= p\left(\mu|x_1, x_2, \ldots, x_n\right) \\
&= \exp\left(-\frac{1}{2}\left(\frac{(\mu - \mu_0)^2}{\tau^2} + \sum_{i=1}^{n}\frac{(x_i - \mu)^2}{\sigma^2}\right)\right) \\
&= \exp\left(-\frac{1}{2}\left(\frac{\mu^2 - 2\mu\mu_0 + \mu_0^2}{\tau^2} + \sum_{i=1}^{n}\frac{\mu^2 - 2\mu x_i + x_i^2}{\sigma^2}\right)\right) \\
&= \exp\left(-\frac{1}{2}\left(\mu^2\Lambda - 2\mu\eta + \left(\frac{\mu_0^2}{\tau^2} + \sum_{i=1}^{n}\frac{x_i^2}{\sigma^2}\right)\right)\right)
\end{aligned}
$$

Whre:

$$
\Lambda = \frac{1}{\tau^2} + \frac{n}{\sigma^2}
$$

$$
\eta = \frac{\mu_0}{\tau^2} + \frac{n\overline{x}}{\sigma^2}
$$

$$
\overline{x} = \frac{1}{n}\sum_{i=1}^{n} x_i
$$

Thus :

$$
\begin{aligned}
\mu_{post} &= \mathbb{E}\left[\mu|D\right] \\
&= \frac{\eta}{\Lambda} \\
&= \frac{\frac{\mu_0}{\tau^2} + \frac{n\overline{x}}{\sigma^2}}{\frac{1}{\tau^2} + \frac{n}{\sigma^2}} \\
&= \frac{\sigma^2\mu_0 + n\tau^2\overline{x}}{\sigma^2 + n\tau^2} \\
&= \frac{\sigma^2}{\sigma^2 + n\tau^2}\mu_0 + \frac{n\tau^2}{\sigma^2 + n\tau^2}\overline{x}
\end{aligned}
$$

And :

$$
\begin{aligned}
\widehat{\Sigma}_{post}^2 &= \mathbb{V}\left[\mu|D\right] \\
&= \frac{1}{\Lambda} \\
&= \frac{\sigma^2\tau^2}{\sigma^2 + n\tau^2}
\end{aligned}
$$

Indeed, the variance decreases in $\frac{1}{n}$.

**With an a priori on $\sigma^2$ but not on** $\mu$   We get $p\left(\sigma^2\right)$ as an Inverse Gamma form.

**With ans a priori on $\mu$ and $\sigma^2$**    Gaussian a priori on $x$ and $\mu$, Inverse Gamma a priori on $\sigma^2$. Please refer to the chapter 9 of the course handout.

### 9.2.8.3   Generalisation of the Beta distribution

Dirichlet is the conjugate of the Multinomial law.

$$p\left(\theta_1, \theta_2, \ldots, \theta_k\right) = \frac{\Gamma\left(\alpha_1 + \alpha_2 + \ldots + \alpha_k\right)}{\Gamma\left(\alpha_1\right)\Gamma\left(\alpha_2\right)\ldots\Gamma\left(\alpha_k\right)}\theta_1^{\alpha_1-1}\theta_2^{\alpha_2-1}\ldots\theta_k^{\alpha_k-1}d\mu\left(\theta\right)$$

Where $\mu$ stands for the uniform measure on $\left\{s \in \mathbb{R}^k \mid \sum_i s_i = 1 \,;\, \forall i,\, s_i \geq 0\right\}$ (simplex).

## 9.3   Model Selection

### 9.3.1   Introduction

Let's consider two models $M_1 \subset M_2$ with $\Theta_1 \subset \Theta_2$. We define:

$$\widehat{\Theta}_{M_i} = \arg \max_{\theta \in \Theta_i} \log \left( p_\theta \left( x_1, x_2, \ldots, x_n \right) \right)$$

where $i \in \{1, 2\}$.

We can't use the maximum likelihood as a score since we have by definition:

$$\log \left( p_{\widehat{\Theta}_{M_2}} \right) \geq \log \left( p_{\widehat{\Theta}_{M_1}} \right)$$

.

We are interested in the capacity of the generalisation of the model: we'd like to avoid over-fitting. Commonly, one way of dealing with that task is to select the size of the model by cross-validation. Here, we'll not develop it furthermore.

In this part we present the *Bayes factors*, which gives us the main Bayes principal for selecting models. Also we will show the link with the penalised version BIC, (Bayesian Information Criterion) which is used by the frequentists so as to "correct" the maximum likelihood and which has good proprieties. The issue with the selection model ask is the issue with the selection of the variables which are an active topic of research. There are others ways of penalising the maximum likelihood and of selecting models.

If $p_0$ is the distribution of the real data, we wish to choose between difference models $(M_i)_{i \in I}$ by maximising $\mathbb{E}_{p_0} \left[ \log \left( p_{M_i} \left( X^* | D \right) \right) \right]$, where $X^*$ is a new test sample distributed as $p_0$ (in fact, it's still the maximum likelihood principle but we take the expectation on new data).

In the Bayesian framework, we can compute the marginal probability of data for a given model

$$\int p \left( x_1, x_2, \ldots, x_n | \theta \right) p \left( \theta | M_i \right) d\theta = p \left( D | M_i \right)$$

and, by applying the Bayes rule, compute the aposteriori probability of the model:

$$p \left( M_i | D \right) = \frac{p \left( D | M_i \right) p \left( M_i \right)}{p \left( D \right)}$$

### 9.3.2   Bayes Factor

Let's introduce the Bayes factors, which enables us to compare two models :

$$\frac{p \left( M_1 | D \right)}{p \left( M_2 | D \right)} = \frac{p \left( D | M_1 \right) p \left( M_1 \right)}{p \left( D | M_2 \right) p \left( M_2 \right)}$$

The marginal probability of data

$$p\left(D|M_i\right) = p\left(x_1, x_2, \ldots, x_n|M_i\right)$$

can decompose itself in a sequential way by using:

$$p\left(x_n|x_1, x_2, \ldots, x_{n-1}, M\right) = \int p\left(x_n|\theta\right) p\left(\theta|x_1, x_2, \ldots, x_{n-1}, M\right) d\theta.$$

Indeed, we get:

$$p(D|M) = p(x_n|x-1, \ldots, x_{n-1}, M)\, p(x_{n-1}|x-1, \ldots, x_{n-2}, M) \ldots p(x_1|M)$$

Such as

$$\frac{1}{n} \log p\left(D|M_i\right) = \frac{1}{n} \sum_{i=1}^{n} \log p(x_i|x_1, \ldots, x_{i-1}, M) \simeq \mathbb{E}_{p_0}\left[\log p_M\left(X|D\right)\right]$$

### 9.3.3   Proposition

the Bayesian score is approximated by the BIC ("Bayesian information criterion").

$$\log p\left(D|M\right) = \log p_{\widehat{\theta}_{MV}}\left(D\right) - \frac{K}{2} \log\left(n\right) + O\left(1\right)$$

With $p_{\widehat{\theta}_{MV}}\left(D\right)$ the data's distribution when the parameter is the maximum likelihood estimator $\widehat{\theta}_{MV}$, $K$ is the number of parameters of the model and $n$ the number of observations.

In the following section, we outline the proof of this result in the case of an exponential family given by $p\left(x|\theta\right) = \exp\left(\langle\theta, \phi\left(X\right)\rangle - A\left(\theta\right)\right)$.

### 9.3.4   Laplace's Method

$$p\left(D|M\right) = \int \prod_{i=1}^{n} p\left(x_i|\theta\right) p\left(\theta\right) d\theta$$

$$= \int \exp\left(\langle\theta, n\bar{\phi}\rangle - n\,A\left(\theta\right)\right) p\left(\theta\right) d\theta$$

$$\langle\theta, n\bar{\phi}\rangle - n\,A(\theta) = \langle\widehat{\theta}, n\bar{\phi}\rangle - n\,A(\widehat{\theta}) + \langle\theta - \widehat{\theta}, n\bar{\phi}\rangle$$

$$- n(\theta - \widehat{\theta})^T \nabla_\theta A(\widehat{\theta}) - \frac{1}{2}(\theta - \widehat{\theta})^T n \nabla_\theta^2 A(\widehat{\theta})(\theta - \widehat{\theta})$$

$$+ \mathrm{R}_n$$

where $\mathrm{R}_n$ is a negligible rest.

But the maximum likelihood is the dual of the maximum entropy : $\max H(p_\theta)$ such that $\mu(\theta) = \bar{\phi}$.

$$\mu(\widehat{\theta}) = \bar{\phi}$$

$$p(D|M) \simeq \exp(\langle \widehat{\theta}, n\bar{\phi} \rangle - n\, A(\widehat{\theta})) \times \int \exp\left(-\frac{1}{2}(\theta - \widehat{\theta})^T n\widehat{\Sigma}(\theta - \widehat{\theta})\right) p(\theta)d\theta$$

However :

1. the information of fisher is equal to $\widehat{\Sigma}^{-1}$

2. $\int \exp\left(-\frac{1}{2}\left(\theta - \widehat{\theta}\right)^T n\widehat{\Sigma}\left(\theta - \widehat{\theta}\right)\right) p(\theta)\, d\theta \simeq c\sqrt{(2\pi)^k \left|\dfrac{\widehat{\Sigma}^{-1}}{n}\right|}$

Thus :

$$\log p\left(D|M\right) = \log p_{\widehat{\theta}}\left(X\right) + \frac{1}{2}\log\left((2\pi)^k \left|\frac{\widehat{\Sigma}^{-1}}{n}\right|\right)$$

$$= \log p_{\widehat{\theta}}\left(X\right) + \frac{k}{2}\log\left(2\pi\right) + \frac{1}{2}\log\left(\left(\frac{1}{n}\right)^k \left|\widehat{\Sigma}^{-1}\right|\right)$$

$$= \log p_{\widehat{\theta}}\left(X\right) + \frac{k}{2}\log\left(2\pi\right) - \frac{k}{2}\log\left(n\right) + \frac{1}{2}\log\left(\left|\widehat{\Sigma}^{-1}\right|\right)$$

The principale reason why presenting the BIC is that a theorem proove the consistancy of the BIC. In other words, when the number of observations is sufficient, thanks to this criterion we choose with a probability that converges to 0, a model that satisfies :

$$M_k \in \mathrm{Argmax}_M\, \mathbb{E}_{p_0}\left[\log\left(p_{\widehat{\theta}_{MV}}\left(X\,;\, M\right)\right)\right]$$