# Bayesian Probabilities

There are two school of thoughts in probability theory —

① Frequentist or Classical → random repeatable events

② Bayesian View → quantification of uncertainity.

There are events which doesn't get repeated numerous times in order to suit the classical or frequentist probability.

Example - how quickly ice cap is melting. In many circumstances, we would like to be able to quantify our expression of uncertainity and make precision precise revisions of uncertainity in presence of new evidence. This can be done using Bayesian interpretation of probability.

Cox in 1946, showed that if numerical values are used to represent degrees of belief, then a simple set of axioms encoding common sense properties of such beliefs leads uniquely to a set of rules for manipulating degrees of belief that are

equivalent to the sum and product rule of probability.

— * —.

In field of pattern recognition. a general notion of probability helps a lot.

Bayes Theorem —

$$p(\omega|D) = \frac{p(D|\omega) \cdot p(\omega)}{p(D)}$$

$\omega$ = parameters

$D$ = observed data = $\{t_1, \dots t_N\}$

This allows us to evaluate the uncertainity in $\omega$ after we have observed $D$ in the form of the posterior probability $p(\omega|D)$.

The quantity $p(D|\omega)$ can be viewed as a function of the parameter vector $\omega$, it is called likelihood function. It expresses how probable the observed data set $b$ for different settings of the parameter vector $\omega$.

$$\text{posterior} \propto \text{likelihood} \times \text{prior}.$$

The denominator is a normalization constant

$$p(D) = \int p(D|w) \cdot p(w) \, dw.$$

In both Bayesian and frequentist paradigms, the likelihood function $p(D|w)$ plays a central role. In frequentist setting, $w$ is considered fixed whose value is determined by some estimator and error bars, considering some distribution on data.

In Bayesian there is only a single dataset $D$ and uncertainity in the parameters is expressed thro' a p.d over $w$.

Maximum likelihood, in which '$w$' is set to the value which maximizes $p(D|w)$. ( negative log likelihood error function, minimization).
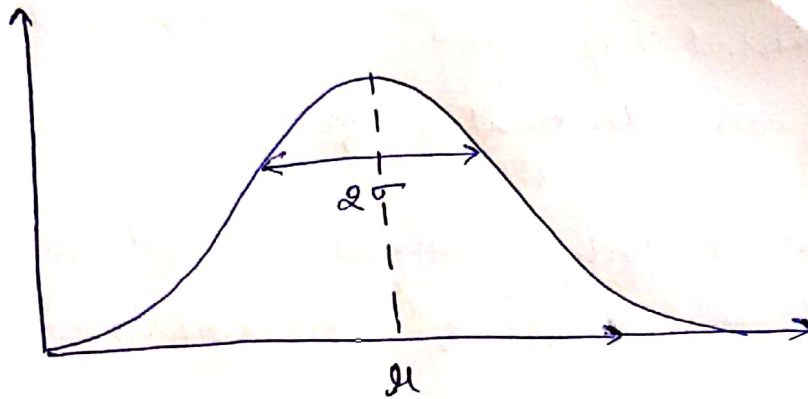
One advantage with Bayesian viewpoint, is inclusion of prior knowledge arises naturally. (talk abud a pair coin toss).

"Monte carlo":

# Gaussian Distribution

$$N(x \mid \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} e^{\left\{ -\frac{(x-\mu)^2}{2\sigma^2} \right\}}$$

$\mu$ = mean , $\sigma$ = std. deviation.



$$\int_{-\infty}^{\infty} N(x \mid \mu, \sigma^2) \, dx = 1.$$

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} N(x \mid \mu, \sigma^2) \, x \, dx = \mu$$

$$\mathbb{E}[x^2] = \int_{-\infty}^{\infty} N(x \mid \mu, \sigma^2) \cdot x^2 \, dx = \mu^2 + \sigma^2$$

$$\text{Var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x] = \sigma^2.$$

The max. of a distribution is known as mode.
For Gaussian, it coincides with mean.

over D- dimensions –

$$N(x|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2}} * \frac{1}{|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right\}$$

D- dimensional vector $\mu \rightarrow$ mean

$\Sigma = D \times D$ dimensional vector $\rightarrow$ covariance.

$|\Sigma|$ determinant of $\Sigma$.

---

Data points that are drawn independently from the same distribution are said to be i.i.d.

i.e independent and identical distribution.

---