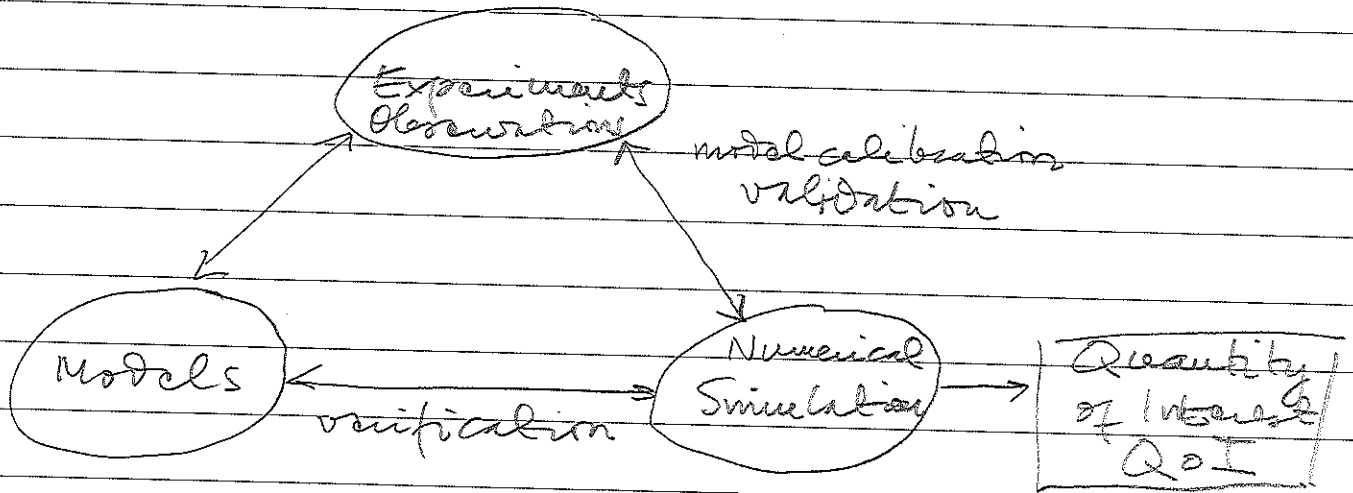


# Predictive science



## Weather prediction

$$\frac{\partial \rho}{\partial t} + \nabla \cdot (\rho \mathbf{v}) = 0$$

$$\frac{\partial \mathbf{v}}{\partial t} = -\mathbf{v} \cdot \nabla \mathbf{v} - \frac{1}{\rho} \nabla p - g \hat{\mathbf{k}} - 2\boldsymbol{\Omega} \times \mathbf{v}$$

$$\rho c_v \frac{\partial T}{\partial t} + \rho \nabla \cdot \mathbf{v} = -\nabla \cdot \mathbf{F} + \nabla \cdot (k \nabla T) + \rho q(T, p, \rho)$$

← concentration of water in solid, liquid, gaseous phase

$$j=1,2,3: \frac{\partial m_j}{\partial t} = -\mathbf{v} \cdot \nabla m_j + S_{m_j}(T, m_j, \chi_j, \rho)$$

← aerosol  $j$  concentration

$$j=1 \dots J: \frac{\partial \chi_j}{\partial t} = -\mathbf{v} \cdot \nabla \chi_j + S_{\chi_j}(T, \chi_j, \rho)$$

phenomenological models for  $S_{m_j}, S_{\chi_j}$

$$S_{m_2} = S_1 + S_2 + S_3 - S_4$$

$$S_1 = \bar{\rho} (m_2 - m_2^*)^2 \left( 1.2 \cdot 10^6 + \left( 1.569 \cdot 10^{-12} \frac{\text{hr}}{d_0 (m_2 - m_2^*)} \right) \right)^{-1}$$

sub-grid modeling

Model and

## Input uncertainties

parameters  $\bar{f}$ ,  $m_2^*$ ,  $v_s$ ,  $d_0$  in  $S_1$   
initial conditions, boundary conditions, model errors

## Numerical errors and uncertainties

grid size: horizontal 5 km, vertical 200 m  
Discretization errors

subgrid models: cloud formation, turbulence

## Measurement errors and uncertainties

limited accuracy of sensors, uncertainty in exact position and time

## Weather forecasts

Determine values, quantify uncertainties  
for initial values, phenomenological parameters  
(Data assimilation, model calibration)

(parameter estimation)

Then run model for forecasts with quantified uncertainties

(multiple)

Ensemble simulations, use initial data and parameters drawn from probability densities

## Definitions

Inputs: parameters, initial conditions, boundary conditions

Quantity of interest (QoI): Output of simulation model (average temperature, precipitation...)

Verification: Quantifying the accuracy of numerical method used to implement math. model

Validation: accuracy of math. model to quantify physical process of interest

## Types of uncertainty

**Alcaotic uncertainty:** statistical, stochastic uncertainty, inherent to a problem, cannot be reduced by more physical or experimental knowledge. Example: initial conditions for weather models

**Epistemic or systematic uncertainty:** caused by simplified models, missing physics, basic lack of knowledge. Example: closure relations, numerical errors

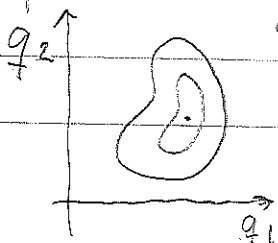
## Predictive estimation

Probabilistic quantification of predicted computational outcome with identified and quantified uncertainties

**Model calibration:** Data to quantify and update uncertainties in parameters, initial and boundary conditions

**Model prediction:** QoI with statistics, probability density function (pdf) for QoI

**Estimation of Validation regime:** contours of constant probability for QoI



$$y = f(X, q), \quad q = (q_1, q_2, \dots, q_p) \text{ para-}$$

random var. (epistemic) model error
indep. var.
metrics

$$Y_i = f(X_i, q) + \delta(X_i) + \epsilon_i \leftarrow \text{random variable (aleaotic)}$$

# Outline of course

	ch.
Typical models	3
Probability, Random processes, statistics	4
Random inputs	5
Parameter selection	6
Frequentist techniques for parameter est.	7
Bayesian techniques for parameter est.	8
Uncertainty propagation	9
Prediction with model discrepancy	12
Surrogate models	13
Local sensitivity analysis	14

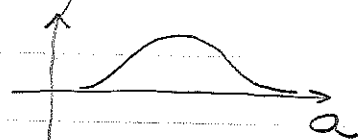
## Simple models

1.  $\frac{dz}{dt} = az + b(t), \quad z(0) = z_0$

parameters  $q = [a, z_0]$

$$z(t, q) = e^{at} \left( z_0 + \int_0^t e^{-as} b(s) ds \right)$$

random variables  $a(\omega), z_0(\omega)$   
 $\omega$  event in probability space



$$\frac{dz}{dt} = a(\omega)z + b(t, \omega), \quad z(0) = z_0(\omega)$$

$$z(t, q) = e^{a(\omega)t} \left[ z_0(\omega) + \int_0^t e^{-a(\omega)s} b(s, \omega) ds \right]$$

## 2. Harmonic oscillator

$$m \frac{d^2 z}{dt^2} + c \frac{dz}{dt} + kz = f_0 \cos \omega_F t$$

$$z(0) = z_0, \quad \frac{dz}{dt}(0) = z_1$$

$$z(t) = c_1 e^{r_1 t} + c_2 e^{r_2 t} + \frac{f_0}{\sqrt{m^2(\omega_0^2 - \omega_F^2)^2 + c^2 \omega_F^2}} \cos(\omega_F t - \delta)$$

$$\cos \delta = \frac{m(\omega_0^2 - \omega_F^2)}{\sqrt{m^2(\omega_0^2 - \omega_F^2)^2 + c^2 \omega_F^2}} \quad \omega_0 = \sqrt{\frac{k}{m}}$$

$$r_{1,2} = \frac{-c \pm \sqrt{c^2 - 4km}}{2m}, \quad z_0, z_1 \Rightarrow c_1, c_2$$

$\text{Re } r_{1,2} < 0 \Rightarrow z(t) \sim \cos \omega_F t - \delta$  for large  $t$   
 If  $c^2 - 4km < 0$  then oscillating first part

$q = [m, c, k, f_0]$  has the same solution as  
 $q = \left[ 1, \frac{c}{m}, \frac{k}{m}, \frac{f_0}{m} \right]$ , only 3 parameters can be uniquely determined by data

measure  $z$  or  $\dot{z}$ , state  $u = [z, \dot{z}]$

observations  $y = C^T u$

either position  $z$   $C^T = (1, 0)$

or velocity  $\dot{z}$   $C^T = (0, 1)$

General form  $y = R(u, q)$   
 description  $\downarrow$  state  $\downarrow$  parameters

all states may not be observable

$$L(q)u = F(q(x)), \quad x \in \Omega = D \times J$$

$x \in D, t \in J$   
 $\mathbb{C}R^2 \quad \mathbb{C}R_1$

$$q(x) = [q_1(x), \dots, q_p(x)]^T \text{ parameters}$$

$$x = [x, t] \text{ independent variables}$$

$$L(q) = [L_1(q), \dots, L_N(q)]^T$$

$$u = [u_1(x), \dots, u_N(x)]^T, \text{ state vector}$$

$$F = [F_1, \dots, F_N]^T \text{ source terms}$$

$$B(q)u = G(q) \quad x \in \partial\Omega$$

boundary and initial conditions

end lecture 1

## Background in probability

### Probability space $(\Omega, \mathcal{F}, P)$

$\Omega$ : sample space, set of all possible outcomes of an experiment

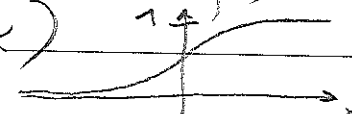
$\mathcal{F}$ :  $\sigma$ -field of subsets of  $\Omega$  containing all events of interest

$P$ :  $\mathcal{F} \rightarrow [0, 1]$ : probability (measure) satisfying

- i:  $P(\emptyset) = 0$   
 ii:  $P(\Omega) = 1$   
 iii: if  $A_i \in \mathcal{F}$  and  $A_i \cap A_j = \emptyset$  then  
 $P(\bigcup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i)$

Realization:  $x$  in  $x = X(\omega)$  for  $\omega \in \Omega$

Cumulative Distribution Function (cdf)

$$F_X(x) = P(\omega \in \Omega \mid X(\omega) \leq x) = P(X \leq x)$$


$$\lim_{x \rightarrow -\infty} F_X(x) = 0, \lim_{x \rightarrow \infty} F_X(x) = 1, x_1 \leq x_2 \Rightarrow F_X(x_1) \leq F_X(x_2)$$

Probability Density Function (pdf)

$$f_X(x), x \in \mathbb{R}$$

$$F_X(x) = \int_{-\infty}^x f_X(s) ds, f_X = \frac{dF_X}{dx}$$

$$f_X(x) \geq 0, \int_{\mathbb{R}} f_X(x) dx = 1$$

$$P(x_1 \leq X \leq x_2) = F_X(x_2) - F_X(x_1) = \int_{x_1}^{x_2} f_X(x) dx$$

Moments  $E(X^n) = \int_{\mathbb{R}} x^n f_X(x) dx$

expected value:  $E(X) = \int x f_X(x) dx = \mu$

second central moment, variance

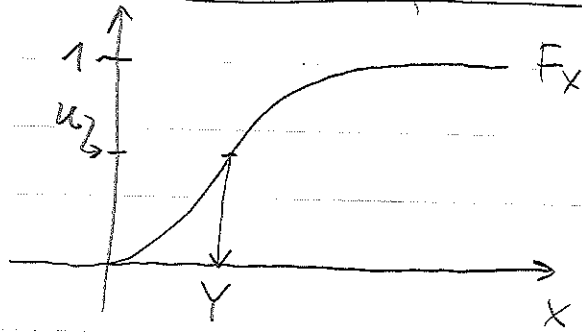
$$\sigma^2 = \text{Var}(X) = E((X - \mu)^2) = \int_{\mathbb{R}} (x - \mu)^2 f_X(x) dx$$

Normal Distribution  $f_X(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-(x-\mu)^2 / 2\sigma^2}$   
 $-\infty < x < \infty$

multivariate normal distribution determinant of  $V$

$$V_{ij} = \text{Cov}(X_i, X_j) \quad f_X(x) = \frac{1}{(2\pi)^n |V|} \exp\left(-\frac{1}{2}(x-\mu)^T V^{-1}(x-\mu)\right)$$

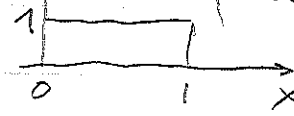
## Inverse Transform Sampling



$X$  has cdf  $F_X(x)$

$U \in \mathcal{U}(0,1)$

uniform distribution



$Y = F_X^{-1}(U)$  has the same distribution as  $X$

$$\begin{aligned} F_Y(y) &= P(Y \leq y) = P(F_X^{-1}(U) \leq y) \\ &= P(U \leq F_X(y)) = F_X(y) \end{aligned}$$

realizations of  $X$ :  $x = F_X^{-1}(u)$

## Multiple random variables

$$X = [X_1, X_2, \dots, X_n]$$

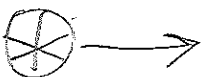
$$\begin{aligned} F_X(x_1, x_2, \dots, x_n) &= P(\omega \in \Omega \mid X_j(\omega) \leq x_j) \\ &= P(X_1 \leq x_1, \dots, X_n \leq x_n) \end{aligned}$$

$$\begin{aligned} \text{cov}(X, Y) &= E[(X - E(X))(Y - E(Y))] \\ &= E(XY) - E(X)E(Y) \end{aligned}$$

$X$  and  $Y$  are independent:  $E(XY) = E(X)E(Y)$   
 $\Rightarrow \text{cov}(X, Y) = 0$

$X_i \sim N(\mu_i, \sigma_i^2)$ , independent

$$Z = \sum_{i=1}^n a_i X_i + b_i \sim N\left(\sum_{i=1}^n (a_i \mu_i + b_i), \sum_{i=1}^n a_i^2 \sigma_i^2\right)$$





Marginal pdf

$X_1, X_2$  have joint pdf  $f_X(x_1, x_2)$

$$f_{X_1}(x_1) = \int_{\mathcal{R}} f_X(x_1, x_2) dx_2$$

Conditional pdf

Conditional density of  $X_1$  given  $X_2 = x_2$  is

$$f_{X_1|X_2}(x_1|x_2) = \begin{cases} \frac{f_X(x_1, x_2)}{f_{X_2}(x_2)} & f_{X_2}(x_2) > 0 \\ 0 & \text{otherwise} \end{cases}$$

Estimator

Parameter  $\theta$ , samples  $X_1, \dots, X_n$

construct estimator of  $\theta$  based on  $X_i$

estimator is random variable

estimate is realization of estimator

unbiased estimator if its mean equals parameter

Examples  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  sample mean

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \text{ sample variance}$$

$$X_i \sim N(\mu, \sigma^2) \Rightarrow \bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

$$S^2 \sim \frac{\sigma^2}{n-1} \chi^2(n-1)$$

$\chi^2$  distribution:  $X_i \sim N(0, 1)$ ;  $X_i$  and  $X_j$  are independent if  $i \neq j$

$\Rightarrow Z = \sum_{i=1}^k X_i^2$  is  $\chi^2(k)$  distributed

$$E(Z) = k, \text{Var}(Z) = 2k$$

### Interval estimator and Confidence interval

$q_L(x), q_R(x)$  such that  $q_L(x) < \eta < q_R(x)$   
(based on  $X = (X_1, \dots, X_n)$ )

$[q_L(x), q_R(x)]$  interval estimator  
+ confidence coefficient  $1 - \alpha$  = confidence interval  
 $P(q_L(X) \leq \eta \leq q_R(X)) = 1 - \alpha$

Interpretation: the frequency of times that the interval will contain parameter  $\eta$

### Model calibration

Statistical model

$Y_i = f(t_i, \theta_0) + \epsilon_i, i = 1, \dots, n$   
random variable, realization  $y_i$

$\epsilon_i$ , errors, independent identically distributed unbiased  $E\{\epsilon_i\} = 0, \text{Var}(\epsilon_i) = \sigma_0^2$  (iid)

### Ordinary Least Squares (OLS) estimator

$$\hat{\theta}_{OLS} = \underset{\theta \in \mathcal{Q}}{\text{argmin}} \sum_{i=1}^n (Y_i - f(t_i, \theta))^2$$

$$\hat{\theta}_{OLS} = \underset{\theta \in \mathcal{Q}}{\text{argmin}} \sum_{i=1}^n (y_i - f(t_i, \theta))^2$$

(ML)

### Maximum Likelihood estimator

$f_Y(y; \theta)$  joint pdf for  $\mathbf{Y} = (Y_1, \dots, Y_n)$   
*(random vector)*

Likelihood function *not conditional pdf, function of  $\theta$*

$$L_Y(\theta) = L(\theta | y) = f_Y(y; \theta)$$

sample  $y$  given,  $\theta$  varies

$n$  iid random variables  $Y_i$

$$L(q|y) = f_Y(y; q) = \prod_{i=1}^n f_Y(y_i; q)$$

log-likelihood function

$$l_Y(q) = \ell(q|y) = \log L(q, y)$$

maximum likelihood estimate (MLE)

$$q_{MLE} = \underset{q \in \mathcal{Q}}{\operatorname{argmax}} \prod_{i=1}^n f_Y(y_i; q)$$

monotonicity of  $\log \Rightarrow$  maximizing  $L(q|y)$  is equivalent to maximizing  $\ell(q|y)$  (computationally easier)

$q_{MLE}$  is the parameter value that makes the output most likely (frequentist view)

### Convergence and Limit Theorems

Sequence  $X_1, X_2, \dots$  converges to  $X$

Almost sure convergence  $X_n \xrightarrow{a.s.} X$

if for every  $\epsilon > 0$

$$\Downarrow \quad P\left(\lim_{n \rightarrow \infty} (X_n - X) < \epsilon\right) = 1$$

Convergence in probability  $X_n \xrightarrow{P} X$

$$\lim_{n \rightarrow \infty} P(|X_n - X| < \epsilon) = 1$$

$\Downarrow$  Convergence in distribution  $X_n \xrightarrow{D} X$

cdf's  $F_{X_1}(x), F_{X_2}(x), \dots, \lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$

$X$  has cdf  $F_X(x)$

## Law of Large numbers

estimator of mean  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$

## Strong Law of Large numbers

$X_i, i=1, 2, \dots$  iid random variables

$$E(X_i) = \mu, \text{Var}(X_i) = \sigma^2 < \infty$$

Then  $P(\lim_{n \rightarrow \infty} |\bar{X}_n - \mu| < \varepsilon) = 1, \bar{X}_n \xrightarrow{\text{a.s.}} \mu$   
(weak form:  $\bar{X}_n \xrightarrow{P} \mu$ )

## Central limit theorem

Assumptions as above

$$Z_n = \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{D} Z \sim N(0, 1)$$

$$\text{or } \bar{X}_n \xrightarrow{D} \bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

## Random processes

temperature  $T$ , heat conductivity  $\alpha(x, \omega)$

$$\frac{\partial T}{\partial x} = \frac{\partial}{\partial x} \left( \alpha(x, \omega) \frac{\partial T}{\partial x} \right),$$

$$\begin{array}{c} T_L \\ | \\ -1 \end{array} \quad \begin{array}{c} | \\ 0 \end{array} \quad \begin{array}{c} T_R \\ | \\ 1 \end{array}, \quad t \geq 0, \quad T(x, 0) = T_0(x)$$

$T(x, t, \omega)$  is a random variable for all  $x, t$   
( $= \{X_t, t \in \mathcal{T}\}$ )

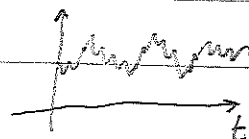
Stochastic process:  $X_t = \{X(t), t \in \mathcal{T}\}$

$\mathcal{T}$  is interval (continuous process)

or consecutive integers (discrete process)

$X_t(\omega)$  realization of process,  $t \in \mathcal{T}, \omega \in \Omega$

sample path or trajectory of  $\omega$



Second order stochastic process:  $E(X_t^2) < \infty$

$\mu(t) = E(X_t), t \in T$

$C(t,s) = Cov(X_t, X_s) = E((X_t - \mu(t))(X_s - \mu(s)))$   
 $t, s \in T$

Markov chains

Discrete state space  $S, X = \{X_i, i \in \mathbb{Z}\}$   
non-unique values

Markov property:  $X_{n+1}$  depends only on  $X_n$

$P(X_{n+1} = x_{n+1} | X_0 = x_0, \dots, X_n = x_n) = P(X_{n+1} = x_{n+1} | X_n = x_n)$

$x_i$  is the state at time  $i$

finite discrete state space

$S = \{x_1, x_2, \dots, x_k\}$ , initial distribution  $p^0$

transition probability  $p_{ij} = P(X_{n+1} = x_j | X_n = x_i)$

$\sum_j p_{ij} = 1$

$P_{ij} = p_{ij}$

$p_j^1 = (p^0 P)_j = \sum_{i=1}^k p_i^0 p_{ij}$

$= \sum_{i=1}^k P(X_1 = x_j | X_0 = x_i) P(X_0 = x_i)$

$= \sum_{i=1}^k P(X_1 = x_j, X_0 = x_i) = P(X_1 = x_j)$

$p^n = p^{n-1} P = p^0 P^n$

(stationary)

Limiting distribution

$\pi = \pi P, \sum \pi_i = 1$

$\pi = \lim_{n \rightarrow \infty} p^n = \lim_{n \rightarrow \infty} p^0 P^n = \lim_{n \rightarrow \infty} p^0 P^{n+1} = (\lim_{n \rightarrow \infty} p^0 P^n) P = \pi P$

⊗ P.T.O.

lit. lect. 2

Random differential equations

Random effects in parameters, initial or boundary conditions

$\frac{dz}{dt} = a(\omega)z + b(t, \omega), z(0) = z_0(\omega)$

Stochastic differential equations

Wieners process

$dZ(t) = -aZ(t)dt + b dW(t)$  or

$Z(t) = Z_0 - \int_0^t aZ(s)ds + \int_0^t b dW(s)$

↳ integral

models Brownian motion

## Statistical inference

Given  $S = \{x_1, x_2, \dots, x_n\}$ ,  $x_j \in \mathbb{R}^N$

$n$  observed realizations of  $X$

want to infer probability distribution of  $X$

parametric inference: prob. distribution has small number of parameters, e.g. mean, variance

problem: estimate parameters

non-parametric inference: construct distribution based on observations only

### Frequentist inference

probabilities are defined as the frequencies with which an event occurs if experiment is repeated a large number of times

OLS or ML estimators

fixed parameter  $\theta_0$ ,  $(1-\alpha) \cdot 100\%$  of confidence intervals contain  $\theta_0$

### Bayesian inference

parameter estimation solution is a posterior probability density

parameters are random variables

credible interval: has  $(1-\alpha) \cdot 100\%$  chance of containing expected parameter

update probability when more data are known

## Bayes' formula

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \frac{P(A, B)}{P(B)}$$

$$\pi(q|y) = \frac{\pi(y|q)\pi_0(q)}{\pi_Y(y)}$$

posterior density      likelihood      prior density      normalization factor

$\pi_0(q)$  quantifies prior knowledge of parameter  $q$   
 no such information: use noninformative prior

e.g.  $\pi_0(q) = \chi_{[0, \infty)}(q)$

$\chi_J$  is indicator function  $\chi_J = \begin{cases} 1, & x \in J \\ 0, & x \notin J \end{cases}$   
 fixed

$\pi(y|q)$  quantifies likelihood  $L(q|y)$  of observing  $y$  given parameter realization  $q$

Joint density  $\pi(q, y) = \pi(y|q)\pi_0(q)$   
 is normalized to 1 by  $\pi_Y(y)$

Posterior density  $\pi(q|y)$  quantifies probability to get parameter  $q$  given observations  $y$

Model calibration with  $\pi_Y(y) = \int_{\mathbb{R}^p} \pi(y|q)\pi_0(q) dq$

$$\Rightarrow \pi(q|y) = \frac{\pi(y|q)\pi_0(q)}{\int \pi(y|q)\pi_0(q) dq}$$

Conjugate priors

prior and posterior distributions have the same parametric form,  $\pi_0(q)$  is then conjugate prior of  $\pi(y|q)$