ECE295, Data Assimilation and Inverse Problems, Spring 2014

Classes

2 April, Intro; Linear discrete Inverse problems (Aster Ch 1) <u>Slides</u>
9 April, SVD (Aster ch 2 and 3) <u>Slides</u>
16 April, Regularization (ch 4) <u>Slides</u>
23 April, Sparse methods (ch 7.2-7.3) <u>Slides</u>
30 April, Bayesian methods (ch 11)
7 May, Markov Chain Monte Carlo (download from <u>Mark Steyvers)</u>
14 May, (Caglar Yardim) Introduction to sequential Bayesian methods
21 May, Kalman
28 May, EKF,UKF,EnKF
4 June, PF

Homework: You can use any programming language, matlab is the most obvious, some Mathematica or Python could be fun!

Hw 1: Download the matlab codes for the book (cd_5.3) from this website <u>http://www.ees.nmt.edu/outside/courses/GEOP529_book.html. Run the 3 examples for chapter 2. Come to class with</u> <u>one question about the examples. Due 9 April.</u>

hw 2, 16 April, SVD ; read chapter one of Steyvers notes,

hw3, 23 April, Regularization

hw4, 30 April, Sparse problems

hw5, 7 May, Monte Carlo integration

hw6, 14 May, Metropolis algorithm (Steyvers Ch 2.4.1-2.4.5); Metropolis-Hastings algorithm, (Steyvers Ch 2.6.1-2.6.3).

Greedy Search Method: Matching Pursuit

Select a column that is most aligned with the current residual



- - Update $S^{(i)}$: If $l \notin S^{(i-1)}, S^{(i)} = S^{(i-1)} \bigcup \{l\}$. Or, keep $S^{(i)}$ the same

• Update
$$r^{(i)}$$
: $r^{(i)} = \mathsf{P}_{a_l}^{\perp} r^{(i-1)} = r^{(i-1)} - a_l a_l^{\top} r^{(i-1)}$

Example 7.4

Consider the recovery of a signal, **m**, shown in Figure 7.18. This 10-s long time series of n = 1001 time points, t_i , is sampled at 100 samples/s and consists of two sine waves at $f_1 = 25$ and $f_2 = 35$ Hz:

$$m_i = h_i \cdot \left(5 \, \cos(2\pi f_1 t_i) + 2 \, \cos(2\pi f_2 t_i) \right) \quad 1 \le i \le n, \tag{7.24}$$

where the signal envelope has also been smoothed with term-by-term multiplication by a Hann taper function,

$$h_i = \frac{1}{2} \left(1 - \cos(2\pi (i-1)/n) \right) \quad 1 \le i \le n.$$
(7.25)





Figure 7.19 Signal recovery using second-order Tikhonov regularization. Solution amplitudes are normalized to improve legibility.



Figure 7.20 A representative solution using second-order Tikhonov regularization that approximately satisfies the discrepancy principle from Figure 7.19 ($\alpha = 10$).

Problem 7.4

$$d_{1000x1} = W_{1000x1000} m_{1000x1} + n_{1000x1}$$
$$q_{100x1} = G_{100x1000} d_{1000x1} = G_{100x1000} W_{1000x1000} m_{1000x1} + n'_{1000x1000} m_{1000x1}$$

-

L1 solutions



Figure 7.21 Signal recovery using compressive sensing with 100 signal measurements. Solution amplitudes are normalized to improve legibility.



Figure 7.22 A representative solution obtained from Figure 7.21 using compressive sensing with $\alpha = 100$ that approximately satisfies the discrepancy principle.



Generating samples from an arbitrary posterior PDF

The rejection method





Generating samples from the posterior PDF

Rejection method



But this requires us to know P_{max}

Step 1: generate a uniform random variable, x_i between a and b

$$p(x_i) = \frac{1}{(b-a)}, \quad a \le x_i \le b$$

Step 2: generate a second uniform random variable, y_i

$$p(y_i) = rac{1}{p_{max}}, \quad 0 \le y_i \le p_{max}$$

Step 3: accept x_i if $y_i \leq p(x_i|d)$ otherwise reject

Step 4: go to step 1

Monte Carlo integration

Consider any integral of the form

$$I = \int_{\mathcal{M}} f(\boldsymbol{m}) p(\boldsymbol{m}|d) d\boldsymbol{m}$$



Given a set of samples m_i (i=,..., N_s) with sampling density $h(m_i)$, the Monte Carlo approximation to I is given by



If the sampling density is proportional to $p(m_i | d)$ then,

$$h(\boldsymbol{m}) = N_s \times p(\boldsymbol{m}|\boldsymbol{d})$$

$$\Rightarrow I \approx \frac{1}{N_s} \sum_{i=1}^{N_s} f(\boldsymbol{m}_i)$$

The variance of the $f(m_i)$ values gives the numerical integration error in I



Finding the area of a circle by throwing darts

$$I = \int_A f(\boldsymbol{m}) d\boldsymbol{m}$$



 $f(\boldsymbol{m}) = \begin{cases} 1 & \boldsymbol{m} \text{ inside circle} \\ 0 & \text{otherwise} \end{cases}$ $h(\boldsymbol{m}) = \frac{N_s}{A}$

$$I \approx \frac{1}{N_s} \sum_{i=1}^{N_s} f(\boldsymbol{m}_i)$$

 \approx Number of points inside the circle

Total number of points

Monte Carlo integration

We have

$$I = \int_{\mathcal{M}} f(\boldsymbol{m}) p(\boldsymbol{m}|d) d\boldsymbol{m} \approx \sum_{i=1}^{N_s} \frac{f(\boldsymbol{m}_i) p(\boldsymbol{m}_i|d)}{h(\boldsymbol{m}_i)} \approx \frac{1}{N_s} \sum_{i=1}^{N_s} f(\boldsymbol{m}_i)$$

The variance in this estimate is given by

$$\sigma_I^2 = rac{1}{N_s} \left\{ rac{1}{N_s^2} \sum_{i=1}^{N_s} f^2(\boldsymbol{m}_i) - \left(rac{1}{N_s} \sum_{i=1}^{N_s} f(\boldsymbol{m}_i)
ight)^2
ight\}$$

- To carry out MC integration of the posterior we ONLY NEED to be able to evaluate the integrand up to a multiplicative constant.
- As the number of samples, N_s, grows the error in the numerical estimate will decrease with the square root of N_s.
 - In principal any sampling density h(m) can be used but the convergence rate will be fastest when $h(m) \propto p(m \mid d)$.

What useful integrals should one calculate using samples distributed according to the posterior p(m | d)?

In low dimensions, these volume and radius formulas simplify to the following:

Dimension	Volume of a ball of radius R	Radius of a ball of volume V
0	1	All balls have volume 1
1	2R	V/2
2	πR^2	$\frac{V^{1/2}}{\sqrt{\pi}}$
3	$\frac{4}{3}\pi R^3$	$\left(\frac{3V}{4\pi}\right)^{1/3}$
4	$\frac{\pi^2}{2}R^4$	$\frac{(2V)^{1/4}}{\sqrt{\pi}}$
5	$\frac{8\pi^2}{15}R^5$	$\left(\frac{15V}{8\pi^2}\right)^{1/5}$
6	$\frac{\pi^3}{6}R^6$	$\frac{(6V)^{1/6}}{\sqrt{\pi}}$
7	$\frac{16\pi^3}{105}R^7$	$\left(\frac{105V}{16\pi^3}\right)^{1/7}$
8	$\frac{\pi^4}{24}R^8$	$\frac{(24V)^{1/8}}{\sqrt{\pi}}$
9	$\frac{32\pi^4}{945}R^9$	$\left(\frac{945V}{32\pi^4}\right)^{1/9}$
10	$\frac{\pi^5}{120}R^{10}$	$\frac{(120V)^{1/10}}{\sqrt{\pi}}$

The volume of a N-dim cube 2^N

For N=2 3.14/2^2=3/4

For N=10 3.14^5/120/2^10=2/1000

This will be hard!



Probabilistic inference

Bayes theorem and all that....







Books











Mass of Saturn

(Laplace 1812)

$$\Pr(x:a\leq x\leq b)=\int_a^b p(x)dx$$

We have already met the concept of using a probability density function p(x) to

describe the state of a random variable.

In the probabilistic (or *Bayesian*) approach, probabilities are also used to describe *inferences* (or *degrees of belief*) about x even if x itself is not a random variable.



Laplace (1812) rediscovered the work of Bayes (1763), and used it to constrain the mass of Saturn. In 150 years the estimate changed by only 0.63% !

But Laplace died in 1827 and then the arguments started...



Bayesian or Frequentist: the arguments





Some thought that using probabilities to describe degrees of belief was too subjective and so they redefined probability as the *long run relative frequency* of a random event. This became the *Frequentist* approach.

To estimate the mass of Saturn the frequentist has to relate the mass to the data through a *statistic*. Since the data contain `random' noise probability theory can be applied to the statistic (which becomes the random variable !). This gave birth to the field of statistics !

But how to choose the statistic ?

.. a plethora of tests and procedures without any clear underlying rationale'

(D. S. Sivia)

`Bayesian is subjective and requires too many guesses' A. Frequentist *`Frequentist is subjective, but BI can solve problems more completely' A. Bayesian*

For a discussion see Sivia (2005, pp 8-11).

Data Analysis: A Bayesian Tutorial' 2nd Ed. D. S. Sivia with J. Skilling, O.U.P. (2005)



Probability theory: Joint probability density functions

A PDF for variable x

p(x)



Probability is proportional to area under the curve or surface

Joint PDF of x and y p(x,y)



If x and y are independent their joint PDF is separable

$$p(x,y) = p(x) \times p(y)$$







Likelihood functions

The likelihood that the data would have occurred for a given model



$$p(d_i|x) = \exp\left\{-\frac{(x - x_{o,i})^2}{2\sigma_i^2}\right\}$$
$$p(d|m) = \exp\left\{-\frac{1}{2}(d - Gm)^T C_D^{-1}(d - Gm)\right\}$$

Maximizing likelihoods is what Frequentists do. It is what we did earlier.

$$\max_{\boldsymbol{m}} p(\boldsymbol{d}|\boldsymbol{m}) = \min_{\boldsymbol{m}} - \ln(p(\boldsymbol{d}|\boldsymbol{m}))$$
$$= \min_{\boldsymbol{m}} (\boldsymbol{d} - \boldsymbol{G}\boldsymbol{m})^T \boldsymbol{C}_D^{-1} (\boldsymbol{d} - \boldsymbol{G}\boldsymbol{m})$$

Maximizing the likelihood = minimizing the data prediction error



Example: Measuring the mass of an object

If we have an object whose mass, *m*, we which to determine. Before we collect any data we believe that its mass is approximately $10.0 \pm 1\mu g$. In probabilistic terms we could represent this as a Gaussian prior distribution Prior PDF

prior
$$p(m) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(m-10.0)^2}$$



Suppose a measurement is taken and a value 11.2μ g is obtained, and the measuring device is believed to give Gaussian errors with mean 0 and $\sigma = 0.5 \mu$ g. Then the likelihood function can be written

$$p(d|m) = \frac{1}{0.5\sqrt{2\pi}} e^{-2(m-11.2)^2}$$
 Likelihood

$$p(m|d) = \frac{1}{\pi} e^{-\frac{1}{2}(m-10.0)^2 - 2(m-11.2)^2} \text{Posterior}$$

$$n(m|d) \propto e^{-\frac{1}{2}(m-10.96)^2}$$

 $p(m|d) \propto e$



The posterior PDF becomes a Gaussian centred at the value of 10.96 μ g with standard deviation $\sigma = (1/5)^{1/2} \approx 0.45$. 259

Example: Measuring the mass of an object

The more accurate new data has changed the estimate of *m* and decreased its uncertainty



For the general linear inverse problem we would have

Prior:
$$p(\boldsymbol{m}) \propto \exp\left\{-\frac{1}{2}(\boldsymbol{m}-\boldsymbol{m}_{o})^{T}C_{m}^{-1}(\boldsymbol{m}-\boldsymbol{m}_{o})\right\}$$

Likelihood: $p(\boldsymbol{d}|\boldsymbol{m}) \propto \exp\left\{-\frac{1}{2}(\boldsymbol{d}-\boldsymbol{G}\boldsymbol{m})^{T}C_{d}^{-1}(\boldsymbol{d}-\boldsymbol{G}\boldsymbol{m})\right\}$
Posterior PDF
 $\propto \exp\left\{-\frac{1}{2}[(\boldsymbol{d}-\boldsymbol{G}\boldsymbol{m})^{T}C_{d}^{-1}(\boldsymbol{d}-\boldsymbol{G}\boldsymbol{m}) + (\boldsymbol{m}-\boldsymbol{m}_{o})^{T}C_{m}^{-1}(\boldsymbol{m}-\boldsymbol{m}_{o})]\right\}$
26



Suppose we have a suspicious coin and we want to know if it is biased or not ?

Let α be the probability that we get a head.

 $\alpha = 1$: means we always get a head. $\alpha = 0$: means we always get a tail. $\alpha = 0.5$: means equal likelihood of head or tail.

We can collect data by tossing the coin many times

 $\{H, T, T, H, \ldots\}$



 $0 \le \alpha \le 1$

We seek a probability density function for α given the data

 $p(\alpha|\boldsymbol{d},I) \propto p(\boldsymbol{d}|lpha,I) imes p(lpha|I)$

Posterior PDF \propto Likelihood x Prior PDF



What is the prior PDF for α ?

Let us assume that it is uniform

$$p(\alpha|I) = 1, \quad 0 \le \alpha \le 1$$



What is the Likelihood function ?

The probability of observing R heads out of N coin tosses is

$$p(\boldsymbol{d}|lpha,I) \propto lpha^R (1-lpha)^{N-R}$$



 $p(\alpha|\boldsymbol{d},I) \propto p(\boldsymbol{d}|lpha,I) imes p(lpha|I)$

Posterior PDF \propto Likelihood x Prior PDF

We have the posterior PDF for α given the data and our prior PDF

$$p(lpha|m{d},I) \propto lpha^R (1-lpha)^{N-R}$$

After N coin tosses let R = number of heads observed. Then we Can plot the probability density for $p(\alpha \mid d)$ as data are collected







But what if three people had different opinions about the coin prior to collecting the data ?

Dr. Blue knows nothing about the coin.

Dr. Green thinks the coin is likely to be almost fair.

Dr. Red thinks the coin is either highly biased to heads or tails.







How to choose a prior ?

An often quoted weakness of Bayesian inversion is the subjectiveness of the prior.

If we know that x > 0 and that $E\{x\} = \mu$. What is an appropriate prior p(x) ?

$$H(X) = \int_{-\infty}^{\infty} p(x) \ln p(x) dx$$
 Entropy

A solution is to choose the prior p(x) that maximizes entropy subject to satisfying the constraints. Using calculus of variations we get

$$p(x) = rac{1}{\mu} e^{-x/\mu}, \quad x \ge 0$$



There is no such thing as a non-informative prior !



Recap: Probabilistic inference

- In the Bayesian treatment, all inferences are expressed in terms of probabilities.
- Bayes' theorem tells us how to combine a priori information with the information provided by the data, and all are expressed as PDFs.
- All Bayesian inference is relative. We always compare what we know after the data are collected to what we know before the data are collected. In practice this means comparing the a posteriori PDF with the a priori PDF.
 - Bayesians argue that this is just a formalization of logical inference.
- Criticisms are that non-informative prior's do not exist, and hence we introduce information if prior's are assumed for convenience.
- The general framework is appealing but can not usually be applied when the number of unknowns is $> = 10^3$.

What can we do with the posterior PDF ?

We could map it out over all of model spaceOnly feasible when dimension of the model space is small.



We seek the maximum of posterior PDF (MAP) and its covariance

$$\phi(\boldsymbol{m}, \boldsymbol{d}) = (\boldsymbol{d} - G \boldsymbol{m})^T C_d^{-1} (\boldsymbol{d} - G \boldsymbol{m}) + (\boldsymbol{m} - \boldsymbol{m}_o)^T C_m^{-1} (\boldsymbol{m} - \boldsymbol{m}_o)$$

 $\max_{\boldsymbol{m}} \phi(\boldsymbol{m}, \boldsymbol{d}); \quad \text{calculate } C_M$

- This is equivalent to the optimization approach earlier.
- Would not make sense when the problem is multi-modal or when the covariance is not representative of its shape.
- We generate model (samples) whose density follows the posterior PDF. Posterior simulation is the main technique used in computational statistics.

In a Bayesian approach the complete Posterior PDF is the answer to the inverse problem, and we always look at its properties.