

- Logistics :
- Attendance
 - 4 homework, mix of written & computational (returned via Gradescope)
 - ↳ (Python, Julia, Matlab).
 - ↳ Stan
 - In-class midterm
 - Final presentation & project, groups of 2 or 3. (+ report, in LaTeX)
 - ↳ overleaf
 - Course website

Other
 - Campuswide?
 -

Overview of course :

1. "Classical numerical analysis"

↳ root finding of function $f(x) = 0$

↳ optimization

↳ think finding

MLE's, or fitting

a NN

$$\left. \begin{matrix} f_1(x_1, \dots, x_n) = 0 \\ \vdots \end{matrix} \right\}$$

↳ numerical linear algebra

↳ solving $A\vec{x} = \vec{b}$

↳ finding eigenvalues.

↳ various factorizations, identities.

↳ function approximation (generally, not interpolation)

- least squares

- regression

- density estimation

19/21

↳ integration

Compute $\int f(x) dx$ numerically

Ex: one might want to compute

marginal distributions $p(x) = \int p(x,y,z) dy dz$

2. statistics, but methods too complicated for pen + paper.

- "simulation" :
 - random variable generation \rightarrow including MCMC
 - estimate confidence intervals
 - regression \rightarrow linear/nonlinear \rightarrow Gaussian Process reg., etc.
 - density estimation
- conditional, hierarchical modelling (large scale Bays).

Example CERN statistical model

Probability Topics Review

see my old webpage
cims.nyu.edu/~rossel/prob20

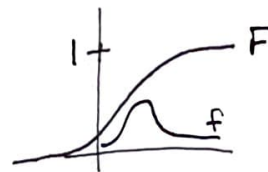
Random Variables: Discrete vs. Continuous

(Refs:
Durrett
or Ross)

less
focus on
these in this
course

each R.V. has a
pdf and distribution
function

$$F(x) = \int_{-\infty}^x f(u) du$$



X is r.v., we say $X \sim \text{Normal}(\mu, \sigma^2)$

$$\Rightarrow f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

for example.

(e.g. exponential,
uniform, Gamma,
etc.
Cauchy)

\Rightarrow Expected value, variance

$$E(X) = \int x f(x) dx$$

$$\text{Var}(X) = E((X - E(X))^2)$$

\Rightarrow Conditional probabilities, independence,

Bayes. (For events)

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

Conditional probability

$$P(A \cap B) = P(A) \cdot P(B)$$

Independence.

$$\Rightarrow P(A \cap B) = P(A|B) P(B)$$

$$= P(B|A) P(A)$$

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

$$= \frac{P(B|A) P(A)}{P(B|A) P(A) + P(B|A^c) P(A^c)}$$

Bayes.

Function of R.V.s

If X has pdf/dist f, F , and

$Y = \varphi(X)$, what is pdf/dist for Y ?

$$\begin{aligned} \Rightarrow G(y) = P(Y \leq y) &= P(\varphi(X) \leq y) = P(X \leq \varphi^{-1}(y)) \\ &= \int_{x \in \varphi^{-1}(y)} f(x) dx \dots \end{aligned}$$

Multiple R.V.s - Joint distributions

X has pdf f ,

Y has pdf g ,

X, Y has joint pdf $h(x, y)$:

$$\Rightarrow P(X \in A, Y \in B) = \int_B \int_A h(x, y) dx dy$$

$$P(X, Y \in \omega) = \iint_{\omega} h(x, y) dx dy$$

Conditional density : $p(x|y) = \frac{h(x, y)}{g(y)}$

$$P(X \in A | Y = y) = \int_A p(x|y) dx$$

Marginal densities : $f(x) = \int h(x, y) dy$

$$g(y) = \int h(x, y) dx.$$

Sums : ~~$Z = X + Y$~~

$$Z = X + Y$$

What is the pdf of Z ?

$$p(z) = \int_x f(z-x) f_Y(x) dx$$

Change of variables

$$u = g(X, Y)$$

$$v = h(X, Y)$$

What is
pdf $f(u, v)$?

Covariance

$$\text{Cov}(X, Y) = E((X - E(X))(Y - E(Y)))$$

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}$$

Multivariate Normal Random Variables — special, important dist.

$$\begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix} \sim N \left(\begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix}, C \right) \quad C_{ij} = \text{Cov}(X_i, X_j)$$

- C is an $n \times n$ matrix, symmetric, positive semi-definite

pdf is given by:

$$f(x_1, \dots, x_n) = \frac{1}{(2\pi)^{n/2} \sqrt{\det C}} e^{-\frac{1}{2}(\vec{x} - \vec{\mu})^T C^{-1}(\vec{x} - \vec{\mu})}$$

linear algebra!

$$\vec{\mu} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix}, \quad \vec{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

If A is an ~~invertible~~ $k \times n$ matrix, then

$\vec{Y} = A\vec{X}$ is also multivariate normal, and

$$E(\vec{Y}) = A\vec{\mu}, \quad \text{and} \quad \text{Cov}(\vec{Y}) = ACAT$$

} normal v.v.
invariant under
linear transformations.

Limit Theorems, inequalities

WLLN:

$$\mu = E(X_i) < \infty$$

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

Then for any $\epsilon > 0$

$$P\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - \mu\right| \geq \epsilon\right) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

(Notation of Convergence)

come back to in
stats review.

CLT

Same assumptions (other assumptions possible)

$$\frac{\frac{1}{n} \sum_{i=1}^n X_i - \mu}{\sigma/\sqrt{n}}$$

$\rightsquigarrow N(0,1)$

↑

convergence in distribution

$$\lim_{n \rightarrow \infty} P\left(\frac{\frac{1}{n} \sum_{i=1}^n X_i - \mu}{\sigma/\sqrt{n}} \leq z\right) = \Phi(z).$$

SLLN

$$\mu = E(X_i) < \infty$$

$$P\left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = \mu\right) = 1.$$

Statistics

Refs: Wasserman, All of Statistics

Rice, Mathematical Statistics & Data Analysis

Casella & Berger, Statistical Inference

Old course website: cims.nyu.edu/~oneil/stat21

- Inference
- Hyp. testing
- Bayesian methods (inference)
- Regression, "curve fitting"
- Density estimation
- Sampling methods (MCMC).

Parametric Methods

Given data x_1, \dots, x_n $\overset{\text{IID}}{\downarrow}$ assumed to be from distribution with pdf $f(x; \theta)$, infer θ . This is inference...

Maximum Likelihood

$L(\theta)$ = joint density evaluated at the data

$$= \prod_i f(x_i; \theta) \leftarrow \text{a function of } \theta.$$

Log-likelihood:

$$\begin{aligned} l(\theta) &= \log L(\theta) \\ &= \sum \log f(x_i; \theta) \end{aligned}$$

$\hat{\theta}$ = MLE is solution to $l'(\theta) = 0$.

How to study the behavior of $\hat{\theta}$?
Assume x_i are random variables...

How to solve $l'(\theta) = 0$ numerically?

Linear Regression



Find line $\hat{y} = \hat{a}x + \hat{b}$ that minimizes "something".

Model: $y = ax + b + \epsilon$
 \uparrow
 IID with $\text{var} \epsilon = \sigma^2$
 $E = 0$

One option: least squares

Minimize $\sum_i (y_i - \hat{y}_i)^2 = \text{RSS}$

- Interpretation:
- ① linear algebra
 - ② statistics (MLE, etc.)
 - ③ calculus

take derivatives, solve optimization problem

solve ~~$A^T a = A^T y$~~ $A^T A \begin{pmatrix} a \\ b \end{pmatrix} = A^T y$

normal equations

when $A = \begin{pmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{pmatrix}$

What is the most accurate way to solve this linear algebra problem?

Other types of "data fitting":

- use basis function ϕ_1, \dots, ϕ_k
- Gaussian process regression
- "smoothing" \leftarrow usually non-parametric

Bayesian Methods

~~Idea~~

Goal: provide density for "belief" on value of a parameter θ
(can get philosophical...)

Idea: Put prior belief on θ , update to posterior belief
by incorporating the data.

Prior: $\theta \sim U(0,1)$

Observe data x_1, \dots, x_n from $f(x; \theta)$

f may also be considered a "prior"

~~Update~~ Update using Bayes Theorem:

$$f(\theta | \vec{x}) = \frac{f(\vec{x} | \theta) f(\theta)}{f(\vec{x})}$$

← prior

← a number

this is the likelihood

$$= \frac{L(\theta) p(\theta)}{C} \propto L(\theta) p(\theta)$$

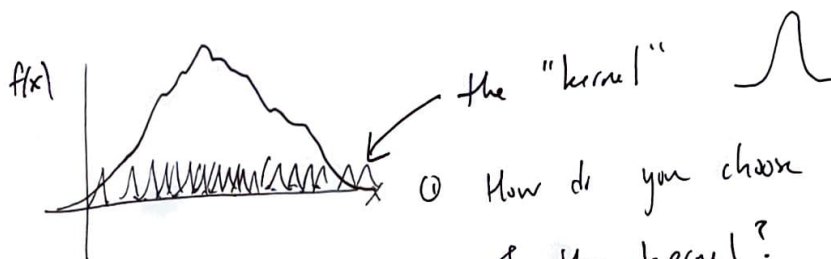
Often $f(\theta | \vec{x})$ is not a known distribution, and
therefore sampling from it to compute, for example,
an expectation (for a point estimate) can be difficult

⇒ the need for computation in statistics

→ mention hierarchical modelling

Density Estimation

Observed data x_1, \dots, x_n , how do we estimate the underlying pdf that generated the data?



① How do you choose the bandwidth of the kernel? (e.g. cross-validation)

② What is a fast algorithm to evaluate $\hat{f}(x)$ at arbitrary x ?

③ Extensions and algorithms for higher dimensions?

④ How do we sample from \hat{f} ?

Sampling Methods

Markov Chain Monte Carlo:

generate a sequence of random variables X_1, X_2, \dots (i.e., a Markov chain) such that as $n \rightarrow \infty$,

$X_n \sim F$, where F is some specified distribution function. Not obvious how to do this in general.