

Statistical Data Analysis - Revision Lec. 1

Exam 9 June 2020 → online, open book
(see email)

No C++ on exam (but, could ask about algorithms)

Format similar to past papers.

Definition & interp. of Probability

↖ Kolmogorov axioms

Interp: . frequency → freq. stat.

. subjective → Bayesian stat.
(degree of belief)

Bayes' thm

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

$$\approx \frac{P(B|A)P(A)}{\sum_i P(B|A_i)P(A_i)}$$

Random variables, prob func., pdfs

pdf $f(x)$ \longrightarrow cdf $F(x) = \int_{-\infty}^x f(x') dx'$

joint pdf $f(x, y)$

marginal "

$$f_x(x) = \int f(x, y) dy$$

conditional pdf

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

Bayes' thm

$$f(x|y) = \frac{f(y|x) f_x(x)}{f_y(y)}$$

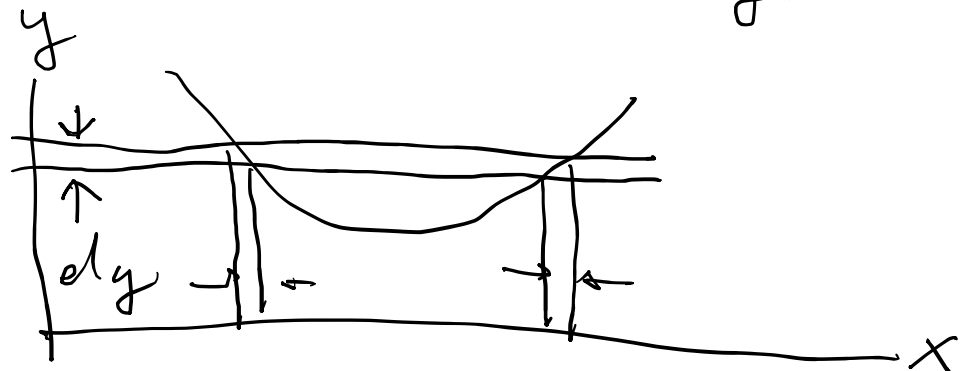
Change of variables

$$x \sim f(x)$$

$y(x)$ is a r.v. $\rightarrow y \sim g(y)$

$$g(y) dy = f(x) dx \Rightarrow g(y) = f(x(y)) \left| \frac{dx}{dy} \right| \quad (1 \rightarrow 1)$$

If $y(x)$ not 1-1



add both
contrib

$$\vec{x} = (x_1, \dots, x_n) \sim f(\vec{x})$$

$$\vec{y} = (y_1(\vec{x}), \dots, y_n(\vec{x})), \quad \text{invert: } \vec{x}(\vec{y})$$

$$g(\vec{y}) = |J| f(\vec{x}(\vec{y}))$$

$$J = \begin{pmatrix} \frac{\partial x_1}{\partial y_1} & \dots & \frac{\partial x_1}{\partial y_n} \\ \vdots & & \vdots \\ \frac{\partial x_n}{\partial y_1} & \dots & \frac{\partial x_n}{\partial y_n} \end{pmatrix}$$

Sums & products of r.v.s.

→ convolutions

$$x, y \sim f(x, y) \quad , \quad z = xy$$

$$g(z) = \int f\left(x, \frac{z}{x}\right) \frac{dx}{x}$$

similar for $z = x + y$

Expectation values (mean)

$$E[x] = \int x f(x) dx \equiv \mu \quad (\text{continuous})$$

$$\text{or } E[n] = \sum_n n P(n) \quad (\text{discrete})$$

$$\text{Variance } V[x] = E[x^2] - (E[x])^2 \equiv \sigma^2$$

$$\text{co } \text{cov}[x, y] = E[xy] - E[x]E[y]$$

mean = centre - of - gravity of pdf

also mode
(pos. of max)



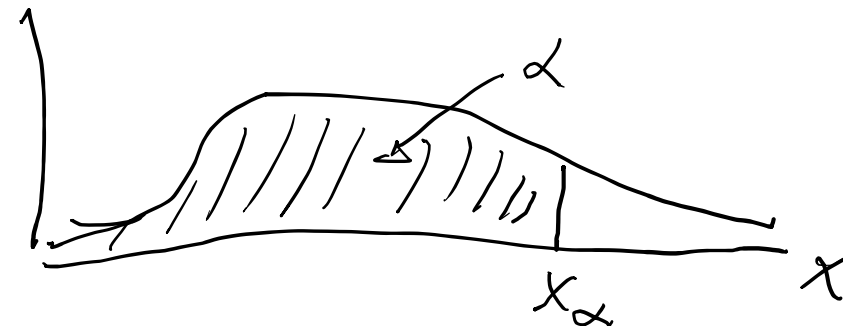
median



α -point

$$F(x_\alpha) = \alpha$$

$$x_\alpha = F^{-1}(\alpha)$$



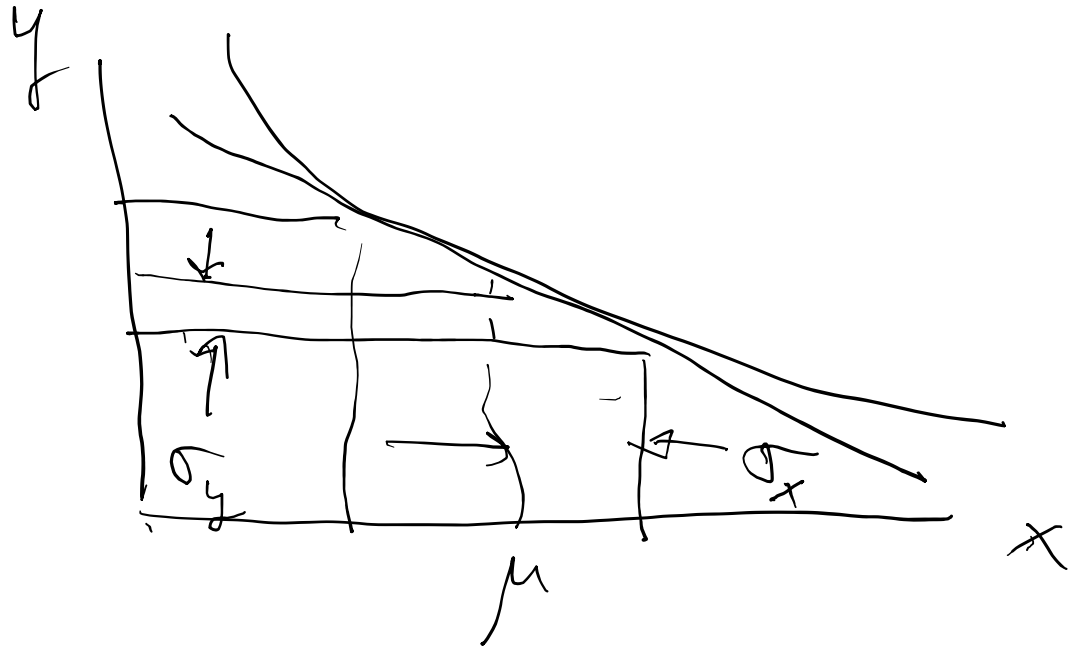
Error propagation

Given \vec{x} , known / estimated $\text{cov}[x_i, x_j] = \chi_{ij}$

\neq " $E[x_i] = \mu_i$

$$V[y(\vec{x})] = \sum_{i,j} \begin{bmatrix} \frac{\partial y}{\partial x_i} & \frac{\partial y}{\partial x_j} \end{bmatrix}_{\vec{x} = \vec{\mu}} \text{cov}[x_i, x_j]$$

Error prop approx, good if $f(\bar{x}) \sim \text{linear}$



Useful pdfs

Gauss $f(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad -\infty < x < \infty$

mean
std. dev.

Expo. $f(x; \lambda) = \frac{1}{\lambda} e^{-x/\lambda} \quad x \geq 0$

mean

Poisson $P(n; \nu) = \frac{\nu^n}{n!} e^{-\nu} \quad n = 0, 1, \dots$

$E[n] = \nu, \quad V[n] = \nu$

Monte Carlo (MC)

1) $r \sim U[0, 1] \rightarrow r_1, r_2, \dots, r_n$
(indep.)

2) $r_1, r_2, \dots, r_n \rightarrow x_1, x_2, \dots, x_n$
such that $x \sim f(x)$

3) Use x_i to estimate property of $f(x)$

e.g. $\int_a^b f(x) dx = \text{fraction of } x_i \text{ in } [a, b]$

Random # gen. (MLCG)

$$n_{i+1} = (an_i) \bmod m$$

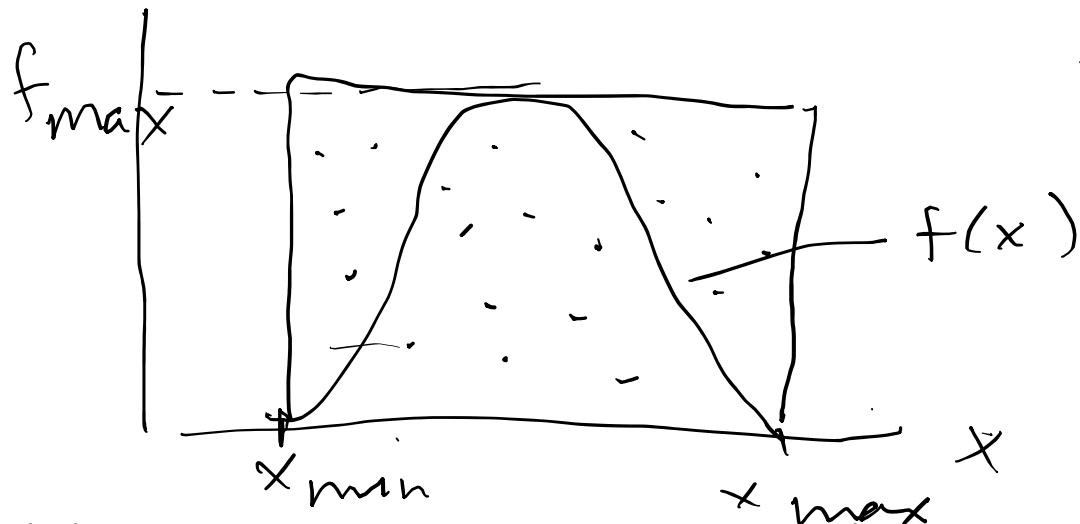
$$r_i = \frac{n_i}{n_{\max}} \sim U[0, 1]$$

↑
 $m-1$

Transformation method $F(x) = r$

∴ solve for $x(r)$

Acceptance - rejection



If point below $f(x)$
accept x .

Statistical tests

Test H_0 ($\leftarrow P(x | H_0)$) vs. H_1

$$P(x \in w | H_0) \leq \alpha \text{ (size)}$$

\uparrow critical region \uparrow small

If x in w ,
reject H_0

$$P(x \in w | H_1) = \text{power of test wrt } H_1 \text{ (} M(H_1) \text{)}$$

H_0 = background evt \leftarrow test H_0

H_1 = signal " "

$$\alpha = P(\text{reject as bkg} | b) = \epsilon_b$$

$$M = P(\text{reject as bkg} | s) = \epsilon_s$$

\uparrow
 \equiv accept as candidate sig.

$$P(s | \text{acc as sig}) = \frac{P(\text{acc as s} | s) \pi_s}{P(\text{acc as s} | s) \pi_s + P(\text{acc as s} | b) \pi_b}$$

Neyman - Pearson Lemma: choose w such that

$$\underbrace{\frac{P(\vec{x} | H_1)}{P(\vec{x} | H_0)}}_{t(\vec{x})} \geq c_\alpha \quad \forall \vec{x} \in w$$
$$\underbrace{\frac{P(\vec{x} | H_1)}{P(\vec{x} | H_0)}}_{t(\vec{x})} < c_\alpha \quad \forall \vec{x} \notin w$$

$t(\vec{x}) = c_\alpha$ boundary of w .

gives max power wrt H_1 , for test of H_0
of size α .

→ multivariate methods
(Machine Learning)

Fisher (linear) $t(\vec{x}) = \sum_i a_i x_i$

Neural Network

Boosted Decision Trees

+ issues (e.g. overtraining)

