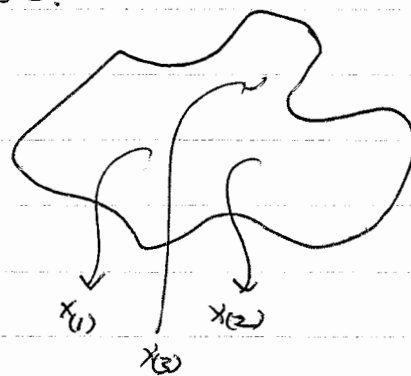


□

Ensembles, Estimators, Gaussians

To model "noise" or variability, we postulate an ensemble Ω of possible observations, + consider each observation as a random draw from Ω .

For univariate quantities, we can characterize Ω by $p(x)$, where $p(x)dx$ is the probability that a random draw is between x + $x+dx$.



Typically we are not interested in determining $p(x)$, but rather "statistics" (functionals) of $p(x)$, e.g.,

$$\text{mean} = \int_{\Omega} x p(x) dx = \langle x \rangle_{\Omega}$$

$$\text{variance} = \int_{\Omega} (x - \langle x \rangle)^2 p(x) dx = \langle (x - \langle x \rangle)^2 \rangle_{\Omega}$$

$$\text{entropy} = \int_{\Omega} -p(x) \ln p(x) dx = - \langle \ln p(x) \rangle_{\Omega}$$

Similarly, for x multivariate ($\vec{x} = (x_1, x_2, \dots, x_N)$).

$$\text{covariance}_{jk} = \int_{\Omega} (x_j - \langle x_j \rangle)(x_k - \langle x_k \rangle) d\vec{x}$$

Time series curve considered as multivariate quantities
 $s(t) = (s(t_1), s(t_2), s(t_3), \dots)$

2

An estimator for a statistic θ is a procedure for taking a set of observations $X_{(1)}, X_{(2)}, \dots, X_{(K)}$ into an estimate θ^{EST} .

We expect that $\theta^{EST}(X_{(1)}, \dots, X_{(K)})$ will not yield the true value $\hat{\theta} = \hat{\theta}(\mathcal{R})$, but it should be "close".

Two kinds of errors

$$\text{BIAS} = \langle \theta^{EST} - \hat{\theta} \rangle_{\text{all draws of } K \text{ samples}}$$

$$\text{VARIANCE} = \langle (\theta^{EST} - \langle \theta^{EST} \rangle_{K \text{ draws}})^2 \rangle_{K \text{ draws}}$$

$$\text{VARIANCE} + (\text{BIAS})^2 = \langle (\theta^{EST} - \hat{\theta})^2 \rangle_{K \text{ draws}}$$

The "plug-in" estimator replaces $\int_{\mathcal{R}}$ by $\frac{1}{K} \sum$.

Plug-in estimator for mean:

$$\theta^{PI} = \frac{1}{K} \sum_{j=1}^K X_{(j)}$$

This is unbiased, since $\langle \theta^{PI} \rangle_{K \text{ draws}} = \frac{1}{K} \sum_{j=1}^K \langle X_{(j)} \rangle_{\mathcal{R}}$

So is the plug-in estimator for any statistic that is linear in \vec{x} , e.g.

$$\sum_r a_r x_r$$

3

The plug-in estimator for the variance is BIA56D

$$\begin{aligned}\Theta^{PI} &= \frac{1}{K} \sum_{j=1}^K \left(X_{(j)} - \frac{1}{K} \sum_{m=1}^K X_{(m)} \right)^2 \\ &= \frac{1}{K} \sum_{j=1}^K X_{(j)}^2 - \frac{1}{K^2} \sum_{j=1}^K \sum_{m=1}^K X_{(j)} X_{(m)} \\ &= \frac{1}{K} \sum_{j=1}^K X_{(j)}^2 - \left(\frac{1}{K} \sum_{j=1}^K X_{(j)} \right)^2\end{aligned}$$

If \mathcal{R} has variance 1, then $\left\langle \sum_{j=1}^K X_{(j)}^2 \right\rangle = K$
 + mean 0

$$\left\langle \left(\sum_{j=1}^K X_{(j)} \right)^2 \right\rangle = K$$

since $\langle X_{(j)} X_{(m)} \rangle = 0$
 on separate draws

$$\begin{aligned}\text{so } \Theta^{PI} &= \frac{1}{K} \cdot K - \frac{1}{K^2} \cdot K \\ &= 1 - \frac{1}{K}, \text{ not } 1.\end{aligned}$$

So we typically use $\frac{\Theta^{PI}}{K-1} = \frac{1}{K-1} \sum_{j=1}^K \left(X_{(j)} - \frac{1}{K} \sum_{m=1}^K X_{(m)} \right)^2$

as an unbiased estimator of the variance.

The above analysis were "lucky" - bias could be calculated independent of \mathcal{R} . But this is not typical. And we also don't know \mathcal{R} (otherwise, we wouldn't need to measure $\vec{x}_{(1)}, \dots$ to estimate Θ).

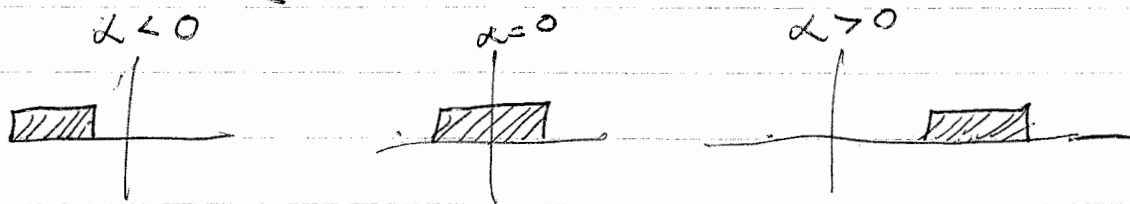
So we need a more general setup:

4

Assume \mathcal{R} is part of a parametric family $\mathcal{R}(\alpha)$.

E.g., [univariate case]


$$\mathcal{R}(\alpha) = \begin{cases} 1, & |x - \alpha| \leq \frac{1}{2} \\ 0, & \text{otherwise} \end{cases}$$



α can be multivariate.

In the above case, the "plug-in" estimator for the mean is unbiased, but has unnecessarily high variance.

A better estimator: choose the highest observation x_{\max}
 " " lowest " x_{\min}
 take $\hat{\Theta} = \frac{1}{2}(x_{\min} + x_{\max})$.

Similarly, if distrib looks like  it is best to discard the extremes before averaging.

Moral: plug-in estimator neither typically unbiased nor least variance

The above set-up leads to 3 generic kinds of estimators (at least):

5

A. Maximum likelihood estimator

Given the observations $\vec{x}_{(1)}, \dots, \vec{x}_{(k)}$, calculate, for each $\vec{\omega}$, the likelihood of this set of observations

$$P(\vec{\omega}) = \prod_{m=1}^K p(\vec{x}_{(m)} \text{ drawn from } \mathcal{N}(\vec{\omega}))$$

Find the $\vec{\omega}_{MLE}$ that maximizes $P(\vec{\omega})$.

$$\text{Then } \Theta^{MLE} = \hat{\Theta}(\mathcal{N}(\vec{\omega}_{MLE})).$$

B. Bayesian estimators

Assume some prior distribution for $\vec{\omega}$, say, $Q(\vec{\omega})$.

Given the observations $\vec{x}_{(1)}, \dots, \vec{x}_{(k)}$, the a priori distribution is modified to an a posteriori distribution

$$Q^{post}(\vec{\omega}) = \frac{P(\vec{\omega}) Q(\vec{\omega})}{\int P(\vec{\omega}') Q(\vec{\omega}') d\vec{\omega}'}$$

$$\text{Can take } \Theta^{BAYES} = \int \hat{\Theta}(\mathcal{N}(\vec{\omega})) Q^{post}(\vec{\omega}) d\vec{\omega}$$

or

$$\Theta^{BAYES, \text{mod}} = \hat{\Theta}(\mathcal{N}(\vec{\beta})) \text{ where } Q^{post}(\vec{\beta}) \text{ is max}$$

$$\text{or } \Theta^{BAYES, \text{gen}} = \frac{\int \hat{\Theta}(\mathcal{N}(\vec{\omega})) g(Q^{post}(\vec{\omega})) d\vec{\omega}}{\int g(Q^{post}(\vec{\omega})) d\vec{\omega}}$$

5

C. "Best" estimator(s) -
use the procedure that minimizes

$$\langle (\theta - \hat{\theta}(x))^2 \rangle_x$$

or, use the procedure that minimizes $\max_x \langle (\theta - \hat{\theta}(x))^2 \rangle$

All of the above depends on choosing $\mathcal{R}(\vec{\theta})$. Usually Gaussians.
Two reasons:

- Central Limit Theorem
- Maximum Entropy property of Gaussians

For a Gaussian, the plug-in estimator for the mean,
and the corrected plug-in estimator for the variance,
are the MLE's.

Bayesian + "best" estimators may differ; $Q(\vec{\theta})$
could be bizarre.

Central Limit Theorem

Informally: Say x_1 is drawn from p_1
 x_2 " " " " p_2
⋮
 x_k " " " " p_k and
no one of these dominates. Then,

□

as $k \rightarrow \infty$, $\frac{1}{k} \sum_{j=1}^k x_j$ is distributed like a

Gaussian.

The plausibility argument will tell us how to make the conditions more precise.

Let $y_2 \in x_1 + x_2$. The p.d.f. for y_2 is

$$q_2(y_2) = \int_{-\infty}^{\infty} p_1(y_2 - x) p_2(x) dx$$

So, with $\hat{p}_1(\omega) = \int e^{-i\omega y} p_1(y) dy$, etc,

$$\hat{q}_2(\omega) = \hat{p}_1(\omega) \hat{p}_2(\omega).$$

Similarly for $y_k = x_1 + \dots + x_k$, $\hat{q}_k(\omega) = \prod_{m=1}^k \hat{p}_m(\omega)$.

Let $z_k = y_k/k$, & say $r_k(z)$ is the p.d.f. for z_k

$$\begin{aligned} \hat{r}(\omega) &= \int e^{-i\omega z} r_k(z) dz \\ &= \int e^{-i\omega z} q_k(kz) \cdot k dz \end{aligned}$$

$$\hat{r}(\omega) = \int e^{-i\omega y/k} q_k(y) dy = \hat{q}_k\left(\frac{\omega}{k}\right) = \prod_{m=1}^k \hat{p}_m\left(\frac{\omega}{k}\right)$$

$$\log \hat{r}(\omega) = \sum_{m=1}^k \log \left(\hat{p}_m\left(\frac{\omega}{k}\right) \right)$$

51

For a p.d.f. $p(x)$, $\log(\hat{p}(w))$ is called the

"characteristic function" $C_p(w)$, & obeys

$$C_{p \times q}(w) = C_p(w) + C_q(w).$$

Write

$$C_p(w) = \sum_{s=0}^{\infty} A_s \frac{(-i)^s}{s!} w^s, \quad (\text{Interpret the } A_s \text{ later})$$

$$\hat{p}(w) = \int e^{-iwx} p(x) dx, \text{ so } \hat{p}(0) = 1$$

$$\hat{p}'(0) = -i \langle X \rangle$$

$$\hat{p}''(0) = -\langle X^2 \rangle$$

$$\frac{d^s}{dw^s} \hat{p}(0) = (-i)^s \langle X^s \rangle.$$

$$\text{So } \hat{p}(w) = \sum_{s=0}^{\infty} \frac{w^s}{s!} (-i)^s \langle X^s \rangle \quad \left[\begin{array}{l} \text{Fine print:} \\ \text{the Taylor series} \\ \text{must exist.} \end{array} \right]$$

$$\hat{p}(w) = 1 - iw \langle X \rangle - \frac{w^2}{2} \langle X^2 \rangle + \frac{iw^3}{6} \langle X^3 \rangle \dots$$

Let's assume that each p_m has mean 0. Then

$$\hat{p}_m\left(\frac{w}{k}\right) = 1 - \frac{w^2}{2k^2} \langle X_m^2 \rangle + \frac{iw^3}{6k^3} \langle X_m^3 \rangle \dots$$

As $k \rightarrow \infty$, since $\log(1+u) = u - \frac{u^2}{2} + \frac{u^3}{3} - \dots$, we can approximate $\log \hat{p}_m\left(\frac{w}{k}\right)$ as:

91

$$\log \hat{p}_m \left(\frac{\omega}{k} \right) = -\frac{\omega^2}{2k^2} \langle X_m^2 \rangle + \frac{i\omega^3}{6k^3} \langle X_m^3 \rangle - \dots$$

$$- \frac{1}{2} \left(-\frac{\omega^2}{2k^2} \langle X_m^2 \rangle + \frac{i\omega^3}{6k^3} \langle X_m^3 \rangle - \dots \right)^2$$

$$+ \frac{1}{3} (\dots)^3 - \dots$$

Only one term that is $O\left(\frac{1}{k^2}\right)$. So

$$\log \hat{r}_m = \sum_{m=1}^k \log \hat{p}_m \left(\frac{\omega}{k} \right) = -\frac{\omega^2}{2k^2} \sum_{m=1}^k \langle X_m^2 \rangle + O\left(\frac{\omega^3}{k^3}\right) + \dots$$

[fine print: need $\sum_{m=1}^k \langle X_m^2 \rangle$ to grow in

proportion to k to be sure that term dominates, similarly for $\langle X_m^3 \rangle$.]

$$* \hat{r}(\omega) \approx e^{-\frac{\omega^2}{2k^2} \sum_{m=1}^k \langle X_m^2 \rangle} \approx e^{-\frac{\omega^2}{2k} V}$$

where V is the average variance of the X_m 's.

$$\text{Can now estimate } r(z) = \frac{1}{2\pi} \int_{-\pi}^{\pi} e^{i\omega z} \hat{r}(\omega) d\omega$$

$$\approx \frac{1}{\sqrt{2\pi}} \frac{1}{\sqrt{V/k}} e^{-\frac{1}{2} z^2 / (V/k)}$$

(we'll do these integrals later)

a Gaussian of variance V/k .

[More fine print: \approx in $*$ means: second moment well approximated, higher moments not well-approximated.]

10

Interpretation of the A_k 's - $A_k(p+q) = A_k(p) + A_k(q)$

$$C_p(\omega) = \log(\hat{p}(\omega)) =$$

"semi-invariants"
"cumulants"

$$-i\omega\langle x \rangle - \frac{\omega^2}{2}\langle x^2 \rangle + \frac{i\omega^3}{6}\langle x^3 \rangle \dots$$

$$- \frac{1}{2} \left(-i\omega\langle x \rangle - \frac{\omega^2}{2}\langle x^2 \rangle + \frac{i\omega^3}{6}\langle x^3 \rangle \dots \right)^2$$

$$+ \frac{1}{3} \left(-i\omega\langle x \rangle - \frac{\omega^2}{2}\langle x^2 \rangle + \frac{i\omega^3}{6}\langle x^3 \rangle \dots \right)^3 \dots$$

$$\sum_{s=0}^{\infty} A_s \frac{(-i)^s}{s!} \omega^s = -i\omega\langle x \rangle - \frac{\omega^2}{2}\langle x^2 \rangle + \frac{1}{2}\omega^2\langle x \rangle^2$$

$$+ \frac{i\omega^3}{6}\langle x^3 \rangle - \frac{1}{2} \left(\frac{i\omega^3}{2} \right) 2\langle x \rangle \langle x^2 \rangle + \frac{1}{3}(-i)^3 \langle x \rangle^3$$

$$+ O(\omega^4)$$

$$\therefore A_1 = -i\langle x \rangle \Rightarrow A_1 = \langle x \rangle$$

$$-\frac{A_2}{2!} = -\frac{1}{2}\langle x^2 \rangle + \frac{1}{2}\langle x \rangle^2 \Rightarrow A_2 = \langle x^2 \rangle - \langle x \rangle^2$$

$$(-i)^3 \frac{A_3}{3!} = \frac{i}{6}\langle x^3 \rangle - \frac{i}{2}\langle x \rangle \langle x^2 \rangle + (i)^3 \frac{\langle x \rangle^3}{3}$$

$$A_3 = \langle x^3 \rangle - 3\langle x \rangle \langle x^2 \rangle + 2\langle x \rangle^3$$

$$A_1 = \text{mean}, \quad A_2 = \text{variance}, \quad \frac{A_3}{(A_2)^{3/2}} = \text{skewness}, \quad \frac{A_4}{(A_2)^2} = \text{kurtosis}$$

III

Some Gaussian integrals

$$I = \int \dots \int e^{-\frac{(\vec{x}-\vec{a})M(\vec{x}-\vec{a})^T}{2}} e^{-\vec{x}\cdot\vec{b}} d\vec{x}$$

where M is a $k \times k$ matrix, symmetric

$$(\vec{x}-\vec{a})M(\vec{x}-\vec{a})^T = \sum_{j,m=1}^k (x_j - a_j)(x_m - a_m) M_{j,m}$$

Say we can write $M = RDR^{-1}$ $\begin{cases} R \text{ real + unitary} \\ D \text{ diagonal,} \\ \text{positive eigenvalues} \end{cases}$
 $= RDR^T$

$$\begin{aligned} (\vec{x}-\vec{a})M(\vec{x}-\vec{a})^T \\ = \vec{y}D\vec{y}^T, \quad \vec{y} = (\vec{x}-\vec{a})R \end{aligned}$$

[if D had an eigenvalue ≤ 0 , then the integral would not $\rightarrow 0$ $\Rightarrow \vec{y} = \alpha \cdot$ corresp eigenvector]

$$\text{Say } D = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_k \end{pmatrix}, \text{ put } z_1 = y_1 \sqrt{\lambda_1}$$

$$\vdots \\ z_k = y_k \sqrt{\lambda_k}$$

$$\vec{z} = \vec{y} \sqrt{D} = (\vec{x}-\vec{a})R\sqrt{D}$$

$$(\vec{x}-\vec{a})M(\vec{x}-\vec{a}) = \vec{z}\vec{z}^T$$

$$\vec{x} = \vec{z} \sqrt{D}^{-1} R^T + \vec{a}$$

$$d\vec{x} = \det(\sqrt{D}^{-1} R) d\vec{z} = \frac{1}{\sqrt{\det M}} d\vec{z} \cdot (\det R = 1)$$

12

$$\begin{aligned}
 \text{So } I &= \int \dots \int e^{-\vec{z}\vec{z}^T/2} e^{(\vec{z}\sqrt{D}^{-1}R^T + \vec{a}) \cdot \vec{b}} \frac{d\vec{z}}{\sqrt{\det M}} \\
 &= \frac{e^{\vec{a} \cdot \vec{b}}}{\sqrt{\det M}} \int \dots \int e^{-\vec{z}\vec{z}^T/2} e^{\vec{z}\sqrt{D}^{-1}R^T \cdot \vec{b}} d\vec{z}
 \end{aligned}$$

Now $\vec{z}\vec{z}^T = z_1^2 + z_2^2 + \dots + z_k^2$

and

$$\vec{z}\sqrt{D}^{-1}R^T \cdot \vec{b} = \sum c_m z_m \text{ for some } \vec{c}^T = \sqrt{D}^{-1}R^T \vec{b}^T$$

$$\text{So } I = \frac{e^{\vec{a} \cdot \vec{b}}}{\sqrt{\det M}} \int \dots \int e^{-(z_1^2 + z_2^2 + \dots + z_k^2)/2} e^{\sum c_m z_m} dz_1 \dots dz_k$$

$$= \frac{e^{\vec{a} \cdot \vec{b}}}{\sqrt{\det M}} \prod_{m=1}^k \int_{-\infty}^{\infty} e^{-\frac{z_m^2}{2} + c_m z_m} dz_m$$

$$= \frac{e^{\vec{a} \cdot \vec{b}}}{\sqrt{\det M}} \prod_{m=1}^k e^{+c_m^2/2} \int_{-\infty}^{\infty} e^{-\frac{1}{2}(z_m^2 - 2c_m z_m + c_m^2)} dz_m$$

$$= \frac{e^{\vec{a} \cdot \vec{b}}}{\sqrt{\det M}} e^{\sum c_m^2/2} (\sqrt{2\pi})^k \left[\int_{-\infty}^{\infty} e^{-u^2/2} du = \sqrt{2\pi} \right]$$

$$\sum c_m^2 = |\vec{c}|^2 = \vec{c} \cdot \vec{c} = (\vec{b} R \sqrt{D}^{-1}) (\sqrt{D}^{-1} R^T \vec{b}^T) = \vec{b} M^{-1} \vec{b}^T$$

$$I = \frac{e^{\vec{a} \cdot \vec{b}}}{\sqrt{\det M}} (2\pi)^{k/2} e^{\vec{b} M^{-1} \vec{b}^T / 2}$$

13

Say $M = V^{-1}$; $\vec{a} = 0$, $\vec{b} = 0$:

$$\int \dots \int e^{-\vec{x} V^{-1} \vec{x}^T / 2} d\vec{x} = \frac{(2\pi)^{k/2}}{\sqrt{\det M}} = \sqrt{\det V} (2\pi)^{k/2}$$

or $\frac{1}{(2\pi)^k} \frac{1}{\sqrt{\det V}} \int \dots \int e^{-\vec{x} V^{-1} \vec{x}^T / 2} d\vec{x} = 1$.

Similar argument* ($V_{km} = V_{mk}$)

$$\frac{1}{(2\pi)^k} \frac{1}{\sqrt{\det V}} \int \dots \int x_k x_m e^{-\vec{x} V^{-1} \vec{x}^T / 2} d\vec{x} = V_{km} \quad (\neq)$$

$$\left[\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} u^2 e^{-u^2/2} du = 1 \right]$$

This establishes

$$\frac{1}{(2\pi)^k} \frac{1}{\sqrt{\det V}} e^{-\vec{x} V^{-1} \vec{x}^T / 2} \text{ as the}$$

p.d.f. of a Gaussian with variances & covariances given by $V = \{V_{km}\}$.

* Similar argument: $\vec{x}^T \vec{x} = (\vec{z} \sqrt{D}^{-1} R^T)^T \vec{z} \sqrt{D}^{-1} R^T$
 $= R \sqrt{D}^{-1} \vec{z}^T \vec{z} \sqrt{D}^{-1} R^T$

$$\text{So } \frac{1}{\sqrt{2\pi}^k} \frac{1}{\sqrt{\det V}} \int \dots \int \vec{x}^T \vec{x} e^{-\vec{x} V^{-1} \vec{x}^T / 2} d\vec{x} =$$

$$R \sqrt{D}^{-1} \left(\frac{1}{\sqrt{2\pi}^k} \frac{1}{\sqrt{\det V}} \int \dots \int \vec{z}^T \vec{z} e^{-\vec{z} \vec{z}^T / 2} \sqrt{\det V} d\vec{z} \right) \sqrt{D}^{-1} R^T$$

$$= R D^{-1} R^T = M^{-1} = V$$

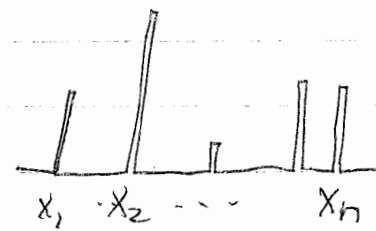
14

$$\text{If } p(\vec{x}) = \frac{1}{(\sqrt{2\pi})^k} \frac{1}{\sqrt{\det V}} e^{-\vec{x} V^{-1} \vec{x}^T / 2},$$

$$\begin{aligned} \hat{p}(\vec{\omega}) &= \int \dots \int p(\vec{x}) e^{-i\vec{\omega} \cdot \vec{x}} d\vec{x} \\ &= e^{-i\omega V \omega^T} \quad [\text{see pg } \boxed{9}] \end{aligned} \quad \left[\begin{array}{l} \vec{a} = 0 \\ M = V^{-1} \\ b = i\omega \\ b^T = -i\omega^T \end{array} \right.$$

Maximum Entropy Distributions

It would be nice to specify $\mathcal{R}(\vec{x})$ with only a few parameters— in particular, would like not to have to specify all moments, on all $p(\vec{x})$'s.

Entropy of a discrete distribution $p(x) =$ 

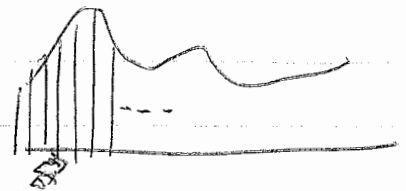
$$p(x) = \sum_{r=1}^n p_r \delta(x - x_r)$$

given by $H(p) = -\sum p_r \log p_r$

If $\log = \log_2$, $H(p) = \#$ of yes-no questions, on average, needed to determine r , if each obs. occurs with frequency p_r .

15

For a continuous distribution:



$$? \quad H(p) \stackrel{?}{=} \lim_{\Delta x \rightarrow 0} \sum_{x_i} (p(x_i) \Delta x) \log p(x_i \Delta x)$$

$$= \lim_{\Delta x \rightarrow 0} \sum_{x_i} p(x_i) \Delta x (\log p(x_i) + \log \Delta x)$$

$$= \lim_{\Delta x \rightarrow 0} \left[\sum_{x_i} p(x_i) \log p(x_i) \Delta x + (\log \Delta x) \sum_{x_i} p(x_i) \Delta x \right]$$

$$= \underbrace{\int p(x) \log p(x) dx}_{\text{"differential entropy" of } p} + \lim_{\Delta x \rightarrow 0} \underbrace{(\log \Delta x)}_{\text{nuisance}}$$

"Nuisance" is irrelevant for comparing entropies.

Plan: Specify only a few moments of p , by \vec{z} .

Then, find $S(\vec{z})$ as the maximum differential entropy ensemble with moments \vec{z} .

16

How to find a maxent distribution, given some constraints?

Discrete case:

$$\text{Entropy} = - \sum p_r \log(p_r)$$

Typical constraint: Mean = $\mu \Rightarrow \sum p_r x_r = \mu$

Variance = $\sigma^2 \Rightarrow \sum p_r (x_r - \mu)^2 = \sigma^2$

The constraints are linear in p .

k^{th} constraint: $\sum_{r=1}^n p_r C_{kr} = A_k$

[We even have a 0^{th} constraint: $\sum_{r=1}^n p_r = 1$, i.e., $C_{0r} = 1$]

Method of Lagrange Multipliers:

Say $f = f(u_1, u_2, \dots, u_n)$ is to be maximized, subject to constraints $c_k(u_1, \dots, u_n) = a_k$.

$$\text{Let } \mathcal{F} = f(u_1, \dots, u_n) + \sum_{k=1}^K \lambda_k c_k(u_1, \dots, u_n).$$

The constrained extremum of \mathcal{F} occurs where $\frac{\partial \mathcal{F}}{\partial u_j} = 0$.

This typically leads to equations $u_j = U_j(\lambda_1, \dots, \lambda_K)$, from which the values of the λ 's must be determined.

17

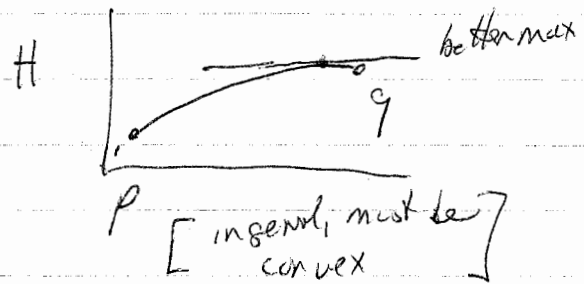
For entropy, the constrained maximum, if it exists, must be unique.

Say p & q are both local maxima that satisfy the constraints, with $H(p) = H(q)$ but $p \neq q$.

$$H((1-m)p + mq) > H(p) \text{ and } H(q)$$

[Mixing property of entropy]

\therefore neither could be the max.



Apply Lagrange Multiplier to $-\sum_{r=1}^n p_r \log p_r$, $\text{c} \sum_{r=1}^n p_r C_{kr}$ as constraints

$$J = -\sum_{r=1}^n p_r \log p_r + \sum_{k=0}^K \lambda_k \sum_{r=1}^n p_r C_{kr}$$

$$\frac{\partial J}{\partial p_j} = -1 + \log p_j + \sum_{k=0}^K \lambda_k C_{kj}$$

$$p_j = B \cdot e^{-\sum_{k=0}^K \lambda_k C_{kj}} \quad \text{Note } C_{0r} = 1$$

So, if we constrain variances & means, the C_{kj} will be:
! (replace variance by 2nd moment)

$$C_{1r} = x_r$$

$$C_{2r} = x_r^2$$

18

Continuous univariate case:

$$\mathcal{F} = - \int p(x) \log p(x) dx + \sum_{k=0}^K \lambda_k \int p(x) C_k(x) dx$$

$$C_{kr} \rightarrow C_k(x) \quad \text{wait } \frac{\partial \mathcal{F}}{\partial p(x)} \quad \left[\text{formally, a "calculus of variations" problem} \right]$$

$$\frac{\partial \mathcal{F}}{\partial p(x)} = -1 - \log p(x) + \sum_{k=0}^K \lambda_k C_k(x)$$

and again

$$p(x) = B e^{\sum_{k=0}^K \lambda_k C_k(x)}$$

Mem: $C_1(x) = x$

Variance: $C_2(x) = x^2$

So $p(x)$ is a Gaussian, - we need to choose λ_1, λ_2 .

Easiest to first shift the problem so that $\mu = 0$, and then use (f) p. 13 to get $\lambda_2 = \frac{1}{2V}$.

Multivariate case: one constraint for each variable's mean
" " " " " variance
" " " " " covariance.

Same strategy.

Conclusion The multivariate Gaussian is the maxent distribution for a constrained mean + variance.