# Practical Statistics for Physicists

Louis Lyons

Oxford &Imperial College

CMS expt at LHC

l.lyons@physics.ox.ac.uk

# Correlations

Basic issue:

For 1 parameter, quote value and uncertainty

For 2 (or more) parameters,

   (e.g. gradient and intercept of straight line fit)

   quote values + uncertainties  **+ correlations**

Just as the concept of variance for single variable is more general than Gaussian distribution, so correlation in more variables does not require multi-dim Gaussian

But more simple to introduce concept this way

# Learning to love the Covariance Matrix

- Introduction via 2-D Gaussian

- Understanding covariance

- Using the covariance matrix

  Combining correlated measurements

- Estimating the covariance matrix

$$y = \frac{1}{\sqrt{2\pi}\,\sigma} \exp\{-(x-\mu)^2/(2\sigma^2)\}$$
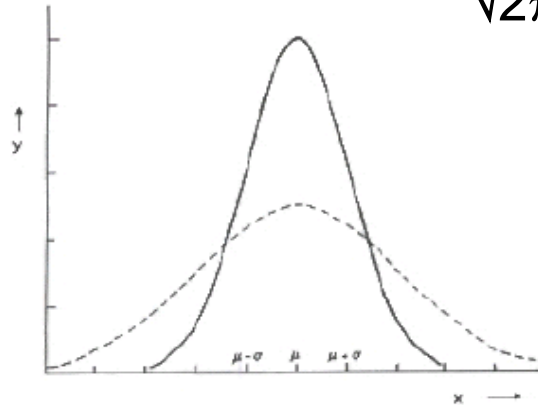
# Reminder of 1-D Gaussian or Normal



Fig. 1.5. The solid curve is the Gaussian distribution of eqn (1.14). The distribution peaks at the mean $\mu$, and its width is characterised by the parameter $\sigma$. The dashed curve is another Gaussian distribution with the same values of $\mu$, but with $\sigma$ twice as large as the solid curve. Because the normalisation condition (1.15) ensures that the area under the curves is the same, the height of the dashed curve is only half that of the solid curve at their maxima. The scale on the $x$-axis refers to the solid curve.

## Significance of σ

i) RMS of Gaussian = σ
(hence factor of 2 in definition of Gaussian)
ii) At $x = \mu \pm \sigma$, $y = y_{max}/\sqrt{e}$ ~0.606 $y_{max}$
(i.e. σ = half-width at 'half'-height)
iii) Fractional area within $\mu \pm \sigma$ = 68%
iv) Height at max = $1/(\sigma\sqrt{2\pi})$

4

## Gaussian in 2-variables

$$P(x) = \frac{1}{\sqrt{2\pi}} \frac{1}{\sigma_x} e^{-\frac{1}{2}\frac{x^2}{\sigma_x^2}}$$

$$P(y) = \frac{1}{\sqrt{2\pi}} \frac{1}{\sigma_y} e^{-\frac{1}{2}\frac{y^2}{\sigma_y^2}}$$

$x + y$ uncorrelated $\Rightarrow$

$$P(x,y) = \frac{1}{2\pi} \frac{1}{\sigma_x \sigma_y} e^{-\frac{1}{2}\left(\frac{x^2}{\sigma_x^2} + \frac{y^2}{\sigma_y^2}\right)}$$

Down on $P(0,0)$ by $e^{-\frac{1}{2}}$ when

$$\frac{x^2}{\sigma_x^2} + \frac{y^2}{\sigma_y^2} = 1$$

Rewrite as

$$(x \quad y)\begin{pmatrix} \frac{1}{\sigma_x^2} & 0 \\ 0 & \frac{1}{\sigma_y^2} \end{pmatrix}\begin{pmatrix} x \\ y \end{pmatrix} = 1$$

Invert $\Rightarrow$ ERROR MATRIX

$$\begin{pmatrix} \sigma_x^2 & 0 \\ 0 & \sigma_y^2 \end{pmatrix}$$

Element $E_{ij}$ - $\langle(x_i - \overline{x_i})(x_j - \overline{x_j})\rangle$

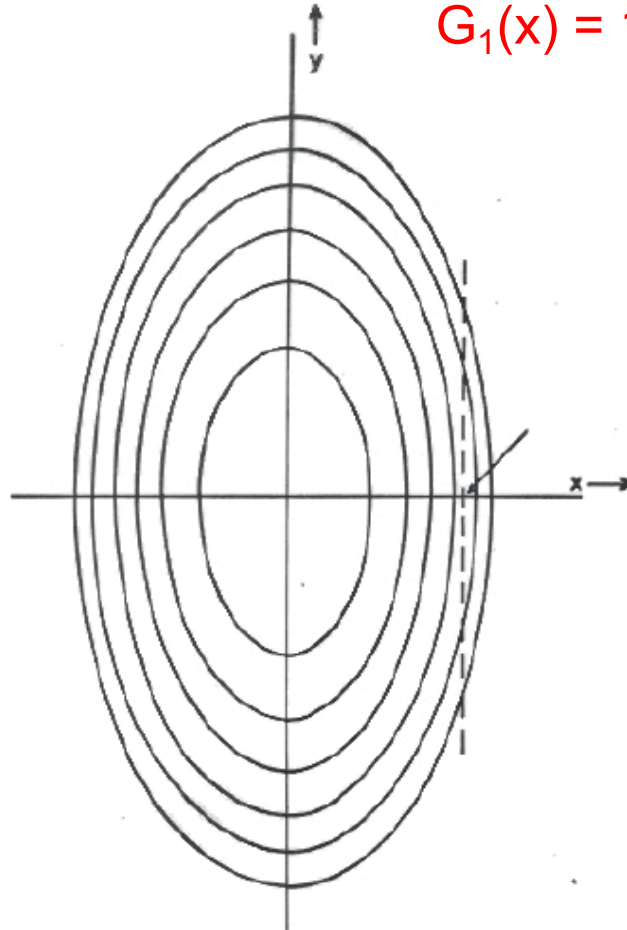Diagonal $E_{ij}$ = variances

Off-diagonal $E_{ij}$ = covariances

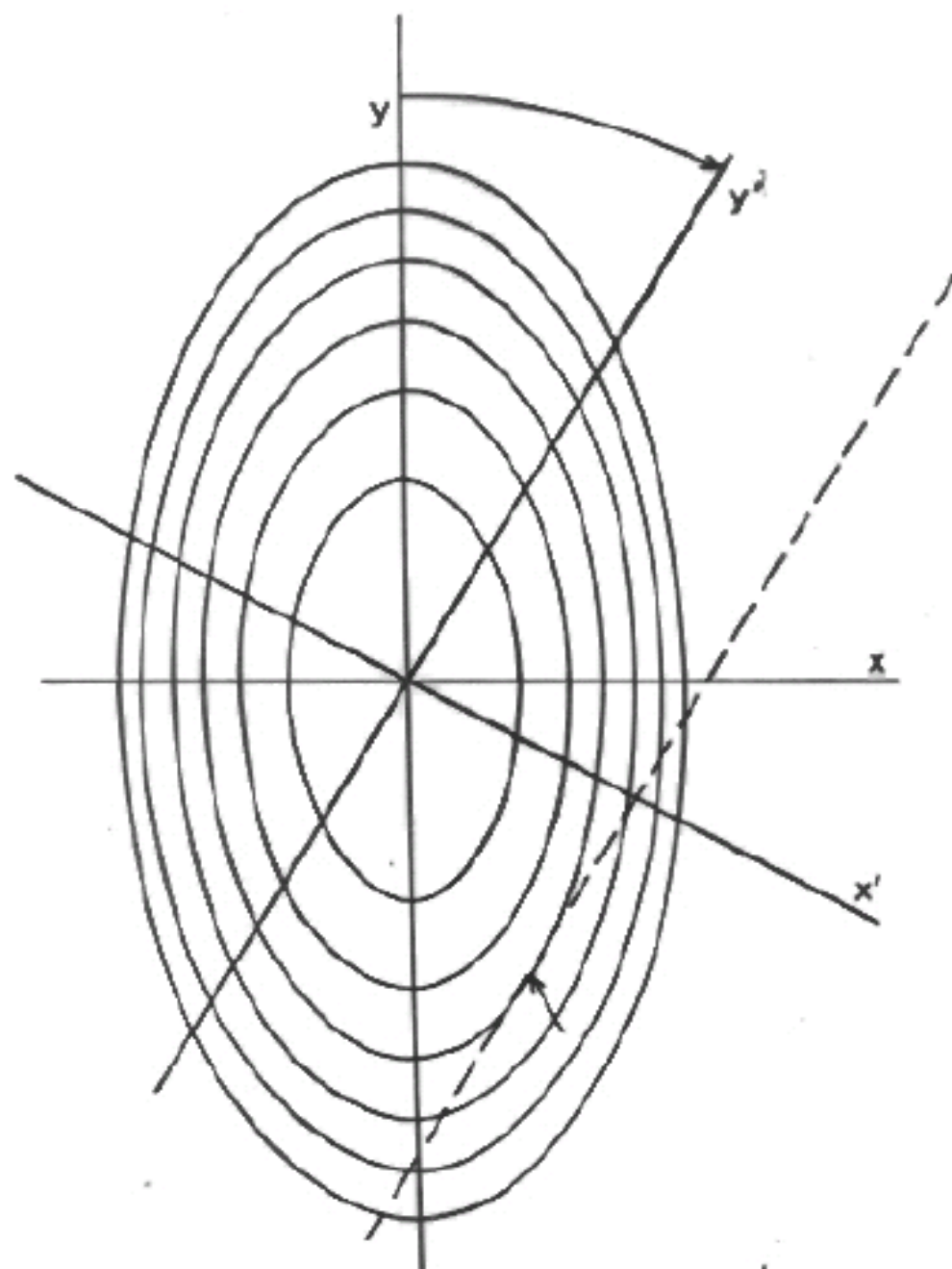# Towards the Covariance Matrix

x and y uncorrelated
$P(x,y\} = G_1(x)\, G_2(y)$
$G_1(x) = 1/(\sqrt{2\pi}\sigma_x)\, \exp\{-x^2/2\sigma_x^2\}$
$G_1(x) = 1/(\sqrt{2\pi}\sigma_y)\, \exp\{-y^2/2\sigma_y^2\}$



$P(x,y) = 1/(2\pi\sigma_x\sigma_y)\, \exp\{-0.5(x^2/\sigma_x^2+y^2/\sigma_y^2)\}$

specific example

$$\sigma_x = \frac{\sqrt{2}}{4} = .354 \qquad\qquad \sigma_y = \frac{\sqrt{2}}{2} = .707$$

Then factor of $e^{-\frac{1}{2}}$ when

$$8x^2 + 2y^2 = 1$$

Now introduce CORRELATIONS by 30° rot$^n$

$$\frac{1}{2}\left[ 13x'^2 + 6\sqrt{3}\,x'y' + 7y'^2 \right] = 1$$

$$\begin{pmatrix} \frac{13}{2} & \frac{3\sqrt{3}}{2} \\[2mm] \frac{3\sqrt{3}}{2} & \frac{7}{2} \end{pmatrix}$$   Inverse Covariance Matrix

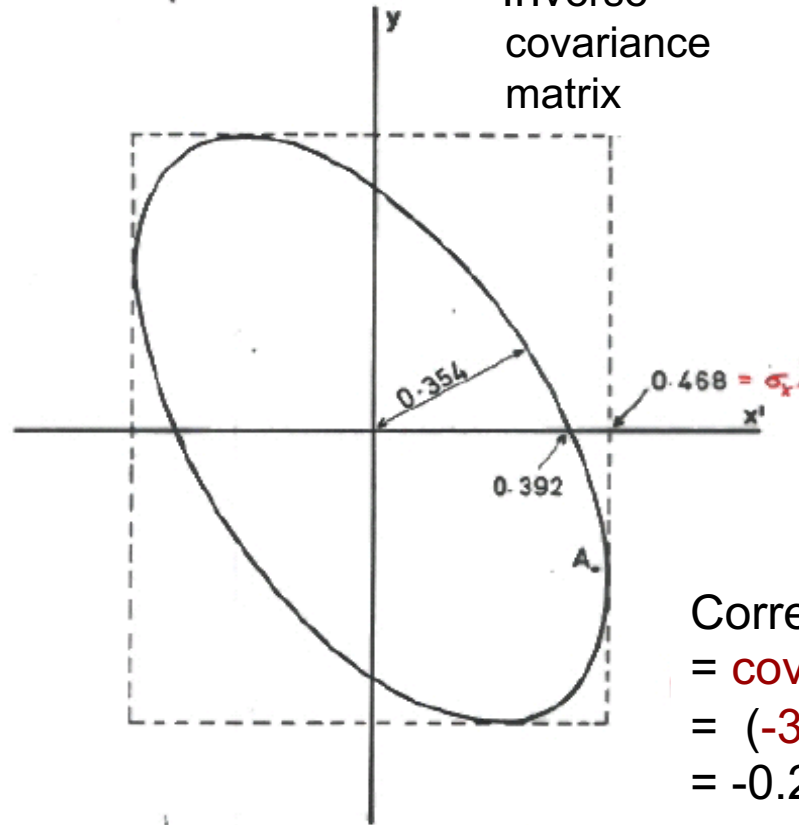$$\frac{1}{32} \times \begin{pmatrix} 7 & -3\sqrt{3} \\ -3\sqrt{3} & 13 \end{pmatrix}$$   Covariance Matrix

8

$$8x^2 + 2y^2 = 1$$
$$0.5(13x'^2 + 6\sqrt{3}x'y' + 7y'^2) = 1$$

$$\begin{pmatrix} 13/2 & 3\sqrt{3}/2 \\ 3\sqrt{3}/2 & 7/2 \end{pmatrix}$$

Inverse
covariance
matrix

$$(1/32)*\begin{pmatrix} 7 & -3\sqrt{3}/2 \\ -3\sqrt{3}/2 & 13 \end{pmatrix}$$

Covariance
matrix



Correlation coefficient $\rho$
= covariance/$\sigma(x')\sigma(y')$
= $(-3\sqrt{3}/2)$/sqrt(7*13)
= -0.27

$7/32 = (0.468)^2 = \sigma(x`)^2$
$1/6.5 = (0.392)^2$
$1/8$ = eigenvalue of covariance matrix = $\sigma(x)^2$

$\left.\begin{array}{c}\sigma_x \\ \sigma_y\end{array}\right\}$ constant

$\rho$ varying

$$\text{Covariance} \begin{pmatrix} \sigma_x^2 & \rho\,\sigma_x\sigma_y \\ \rho\,\sigma_x\,\sigma_y & \sigma_y^2 \end{pmatrix}$$

Covariance matrix,
ρ in range -1→+1



$\rho = 0$

$\rho = -0.9$

$\rho = +1$

# Using the Covariance Matrix

(i)  Function of variables
$$y = y(x_a, x_b)$$
Given covariance matrix for $x_a$, $x_b$, what is $\sigma_y$ ?

Differentiate, square, average

$$\overline{\delta y^2} = \left(\frac{\partial y}{\partial x_a}\right)^2 \overline{\delta x_a^2} + \left(\frac{\partial y}{\partial x_b}\right)^2 \overline{\delta x_b^2} + 2 \frac{\partial y}{\partial x_a}\frac{\partial y}{\partial x_b}\overline{\delta x_a \delta x_b}$$

Zero if $x_a$, $x_b$ uncorrelated

OR

$$\overline{\delta y^2} = \left(\frac{\partial y}{\partial x_a} \quad \frac{\partial y}{\partial x_b}\right)\left(\begin{array}{cc} \overline{\delta x_a^2} & \overline{\delta x_a \delta x_b} \\ \overline{\delta x_a \delta x_b} & \overline{\delta x_b^2} \end{array}\right)\left(\begin{array}{c} \frac{\partial y}{\partial x_a} \\ \frac{\partial y}{\partial x_b} \end{array}\right)$$

$\tilde{D}$

Error matrix

Derivative vector $D$

$$\sigma_y^2 = \tilde{D}ED$$

(ii) Change of variables $x_a = x_a(p_i, p_j)$
$$x_b = x_b(p_i, p_j)$$

e.g Cartesian to polars;   or
Points in x.y → intercept and gradient of line

Given cov matrix for $p_i, p_j$, what is cov matrix for $x_a, x_b$ ?
Differentiate, calculate $\delta x_a \delta x_b$, and average

$$\delta x_a = \frac{\partial x_a}{\partial p_i}\delta p_i + \frac{\partial x_a}{\partial p_j}\delta p_j \qquad (+ \text{ sim for } x_b)$$

Then
$$\overline{\delta x_a^2} = \left(\frac{\partial x_a}{\partial p_i}\right)^2\overline{\delta p_i^2} + \left(\frac{\partial x_a}{\partial p_j}\right)^2\overline{\delta p_j^2} + 2\frac{\partial x_a}{\partial p_i}\frac{\partial x_a}{\partial p_j}\overline{\delta p_i \delta p_j}$$

$$\overline{\delta x_a \delta x_b} = \frac{\partial x_a}{\partial p_i}\frac{\partial x_b}{\partial p_i}\overline{\delta p_i^2} + \frac{\partial x_a}{\partial p_j}\frac{\partial x_b}{\partial p_j}\overline{\delta p_j^2} + \left(\frac{\partial x_a}{\partial p_i}\frac{\partial x_b}{\partial p_j} + \frac{\partial x_a}{\partial p_j}\frac{\partial x_b}{\partial p_i}\right) \times \overline{\delta p_i \delta p_j}$$

$$+ \overline{\delta x_b^2} \quad \text{like} \quad \overline{\delta x_a^2}$$

N.B. Change of variables does not have to be N→N
    e.g. straight line fit involves   N→2
Then i) & ii) are both examples of N→M (M≤N)
       where  M=1  in i)      M=N  in ii)

12

i.e.

$$\begin{pmatrix} \overline{\delta x_a^2} & \overline{\delta x_a \delta x_b} \\ \overline{\delta x_a \delta x_b} & \overline{\delta x_b^2} \end{pmatrix} = \begin{pmatrix} \frac{\partial x_a}{\partial p_i} & \frac{\partial x_a}{\partial p_j} \\ \frac{\partial x_b}{\pa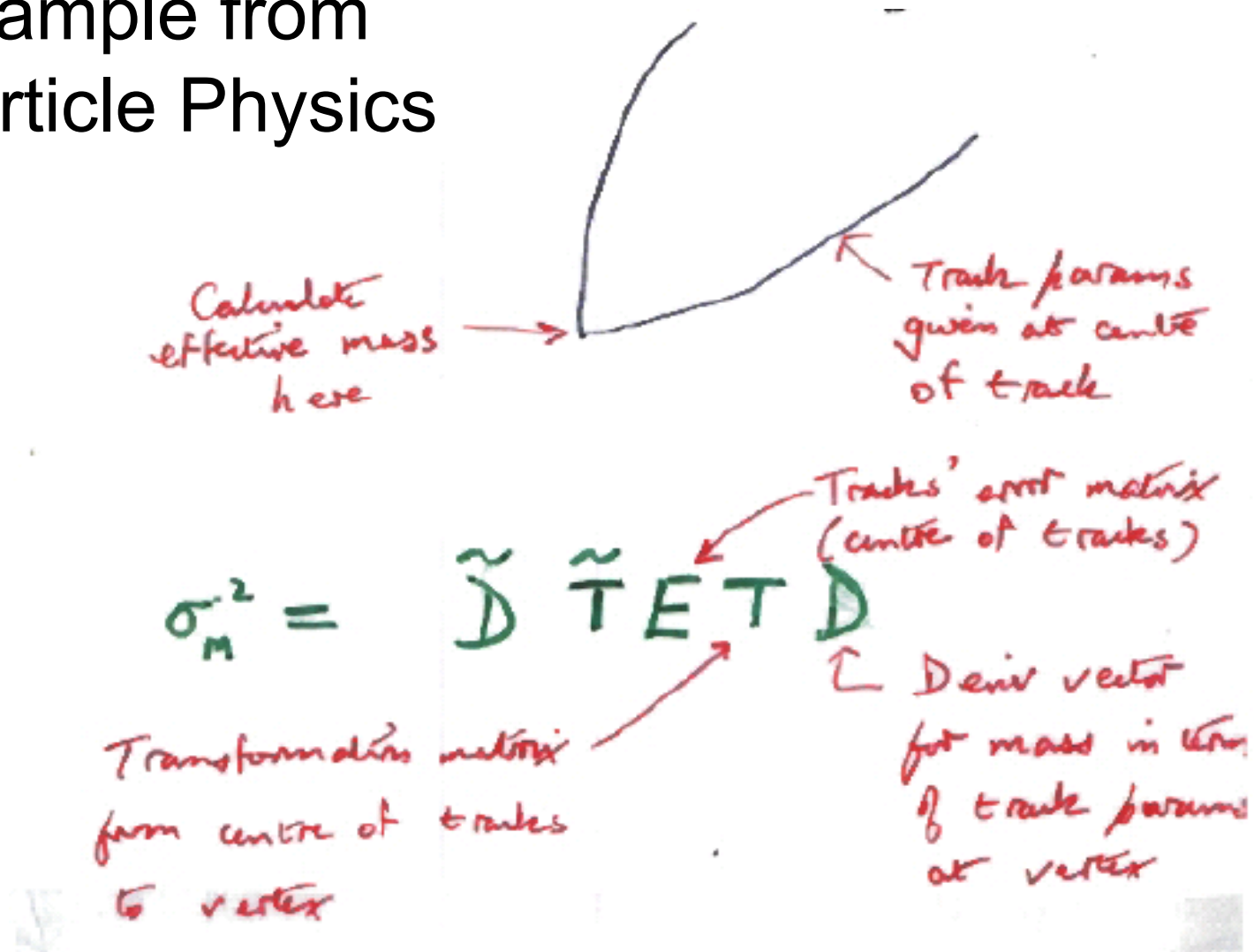rtial p_i} & \frac{\partial x_b}{\partial p_j} \end{pmatrix} \begin{pmatrix} \overline{\delta p_i^2} & \overline{\delta p_i \delta p_j} \\ \overline{\delta p_i \delta p_j} & \overline{\delta p_j^2} \end{pmatrix} \begin{pmatrix} \frac{\partial x_a}{\partial p_i} & \frac{\partial x_b}{\partial p_i} \\ \frac{\partial x_a}{\partial p_j} & \frac{\partial x_b}{\partial p_j} \end{pmatrix}$$

↑ New error matrix

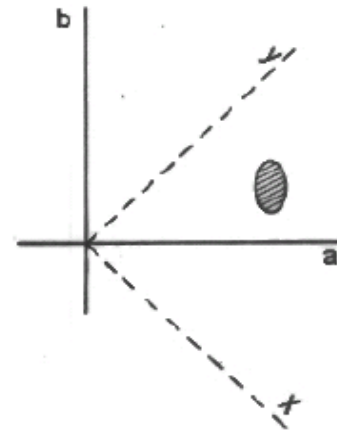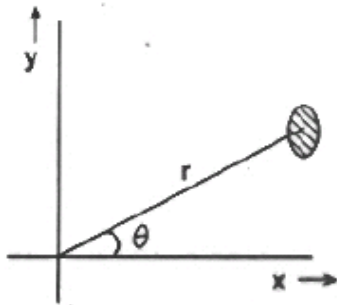↑ $\tilde{T}$

↑ Old error matrix

↑ Transform matrix T

$$E_x = \tilde{T} E_p T$$

**BEWARE!**

# Example from Particle Physics



Calculate effective mass here →

Track params given at centre of track

Tracks' error matrix (centre of tracks)

$$\sigma_M^2 = \tilde{D}\ \tilde{T}\ E\ T\ D$$

Transformation matrix from centre of tracks to vertex

Deriv vector for mass in terms of track params at vertex

14

# Examples of correlated variables

# Using the Covariance Matrix

COMBINING RESULTS

If $a_i \pm \sigma_i$ are independent:

$$\text{Minimise} \quad S = \sum \left( \frac{a_i - \hat{a}}{\sigma_i} \right)^2$$

$$\Rightarrow \quad \hat{a} = \frac{\sum a_i w_i}{\sum w_i} \qquad w_i = 1/\sigma_i^2$$

Now $a_i \pm \sigma_i$ are correlated with error matrix $\underline{\underline{E}}$

$$\underline{\underline{E}} = \begin{pmatrix} \sigma_1^2 & cov(1,2) & cov(1,3) & \cdots \\ cov(1,2) & \sigma_2^2 & cov(2,3) & -- \\ & & \cdots & \cdots & -- & - & \cdot \end{pmatrix}$$

$$S = \sum_{i,j} (a_i - \hat{a}) \, \underline{\underline{E}}^{-1}_{ij} \, (a_j - \hat{a})$$

$\uparrow$ INVERSE ERROR MATRIX

N.B. $\hat{a}$ CAN LIE OUTSIDE $a_i$

$$\sigma_a \rightarrow 0 \quad \text{AS} \quad \rho \rightarrow \pm 1$$

$$\underline{\underline{E}}^{-1} = \begin{pmatrix} 1/\sigma_1^2 & 0 & 0 & \cdots \\ 0 & 1/\sigma_2^2 & 0 & \\ \vdots & \vdots & \vdots & \end{pmatrix} \quad \text{FOR UNCORRELATED}$$
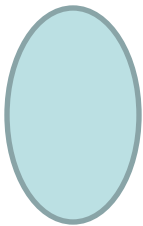
# BLUE
# Best Linear Unbiassed Estimate

Combine several possibly correlated estimates of same quantity
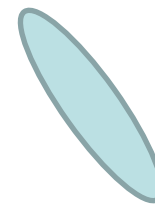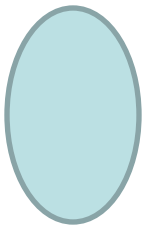e.g. $v_1$, $v_2$, $v_3$

Covariance matrix

$$\begin{pmatrix} \sigma_1^2 & cov_{12} & cov_{13} \\ cov_{12} & \sigma_2^2 & cov_{23} \\ cov_{13} & cov_{23} & \sigma_3^2 \end{pmatrix}$$

Uncorrelated          Positive correlation          Negative correlation

$$cov_{ij} = \rho_{ij}\, \sigma_i\, \sigma_j \quad \text{with} \quad -1 \leq \rho \leq 1$$

# BLUE
# Best Linear Unbiassed Estimate

Combine several possibly correlated estimates of same quantity
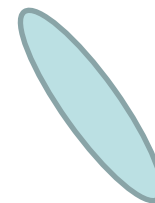e.g. $v_1$, $v_2$, $v_3$

Covariance matrix

$$\begin{pmatrix} \sigma_1^2 & \text{cov}_{12} & \text{cov}_{13} \\ \text{cov}_{12} & \sigma_2^2 & \text{cov}_{23} \\ \text{cov}_{13} & \text{cov}_{23} & \sigma_3^2 \end{pmatrix}$$

Uncorrelated          Positive correlation          Negative correlation

$$\text{cov}_{ij} = \rho_{ij}\, \sigma_i\, \sigma_j \quad \text{with} \quad -1 \le \rho \le 1$$

Lyons, Gibault + Clifford
NIM A270 (1988) 42

$v_{best} = w_1 v_1 + w_2 v_2 + w_3 v_3$ **L**inear

with $w_1 + w_2 + w_3 = 1$ **U**nbiassed

to give $\sigma_{best}$ = min (wrt $w_1, w_2, w_3$) **B**est

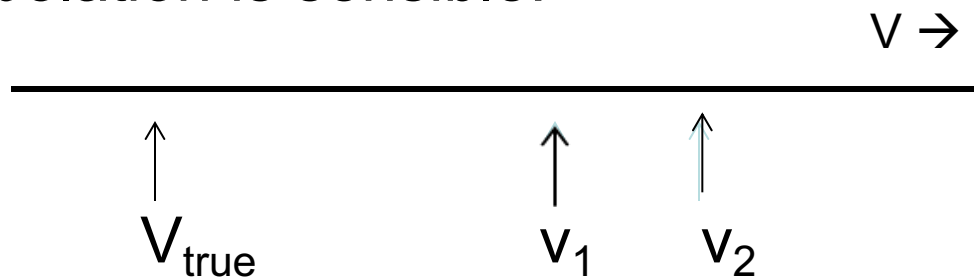For uncorrelated case, $w_i \sim 1/\sigma_i^2$

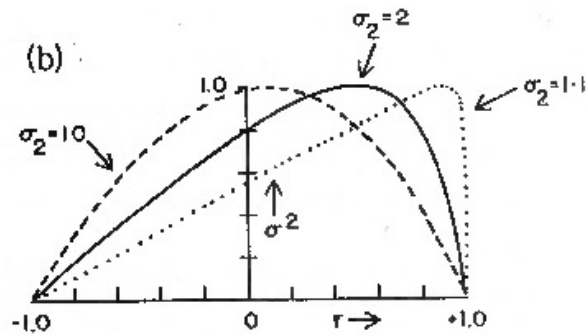For correlated pair of measurements with $\sigma_1 < \sigma_2$

$v_{best} = \alpha v_1 + \beta v_2$ $\qquad \beta = 1 - \alpha$

$\beta = 0$ for $\rho = \sigma_1/\sigma_2$ $\qquad$ (Smaller $\beta \rightarrow$ weights both >0)

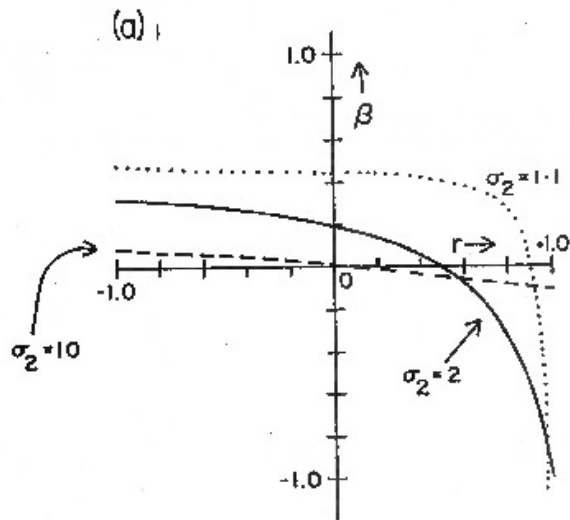$\beta < 0$ for $\rho > \sigma_1/\sigma_2$ i.e. extrapolation! $\qquad$ e.g. $v_{best} = 2v_1 - v_2$

Extrapolation is sensible:
$\qquad\qquad\qquad\qquad\qquad\qquad v \rightarrow$

$V_{true} \qquad\qquad\qquad v_1 \quad v_2$

(b)

Beware extrapolations because

[b] $\sigma_{best}$ tends to zero, for $\rho$ = +1 or -1

(a)

[a] $v_{best}$ sensitive to $\rho$ and $\sigma_1/\sigma_2$

N.B. For different analyses of ~ same data, $\rho$ ~ 1, so choose 'better' analysis, rather than combining

Fig. 1

20

N.B. $\sigma_{best}$ depends on $\sigma_1$, $\sigma_2$ and $\rho$, but not on $v_1 - v_2$

 e.g. Combining  0±3 and x±3  gives x/2 ± 2

<span style="color:blue">BLUE</span> = $\chi^2$

$S(v_{best}) = \Sigma \ (v_i - v_{best}) \ E^{-1}_{ij} \ (v_j - v_{best})$ , and minimise S wrt $v_{best}$

$S_{min}$ distributed like $\chi^2$, so measures Goodness of Fit

But <span style="color:blue">BLUE</span> gives weights for each $v_i$

Can be used to see contributions to $\sigma_{best}$ from each source
of uncertainties e.g. statistical and systematics
                    different systematics

<span style="color:green">Extended by Valassi to combining more than one measured
quantity e.g. intercepts and gradients of a straight  line</span>

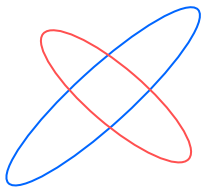# MORE COMBINING:
## Several pairs of correlated params

$$(x_i, y_i) \quad \text{with} \quad \underline{\underline{E}}_i = \begin{pmatrix} \sigma_x^2 & cov \\ cov & \sigma_y^2 \end{pmatrix}_i$$

$$j = \sum_i \left\{ (x_i - \hat{x})^2 E_{11,i}^{-1} + (y_i - \hat{y})^2 E_{22,i}^{-1} + 2(x_i - \hat{x})(y_i - \hat{y}) E_{12,i}^{-1} \right\}$$
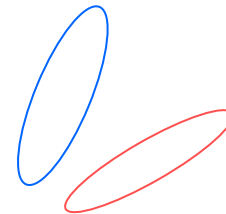
i.e. result: —

Inverse error matrix on result $\hat{x}, \hat{y}$

$$= \sum_i \underline{\underline{E}}_i^{-1}$$

cf $\frac{1}{\sigma^2} = \sum \frac{1}{\sigma_i^2}$ for single uncorrelated meas.
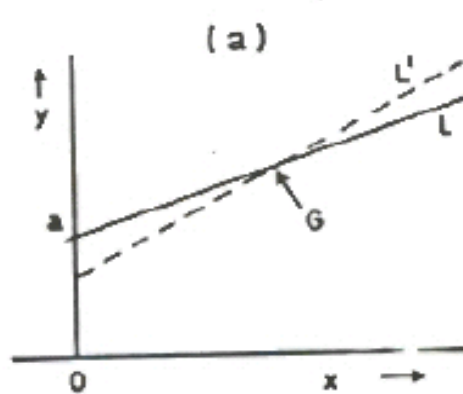


Small uncertainty

Example: Straight line fitting
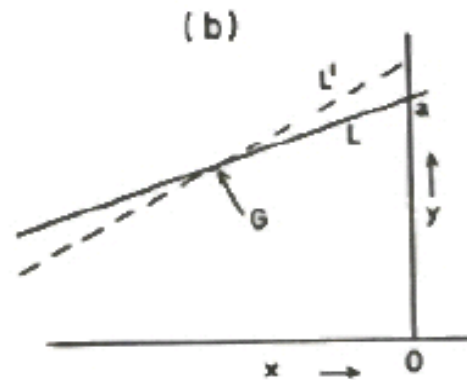


$x_{best}$ outside $x_1 \rightarrow x_2$
$y_{best}$ outside $y_1 \rightarrow y_2$

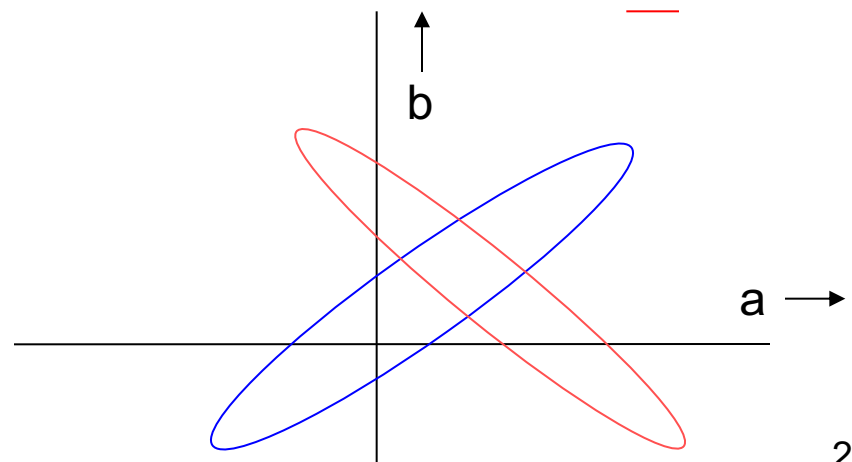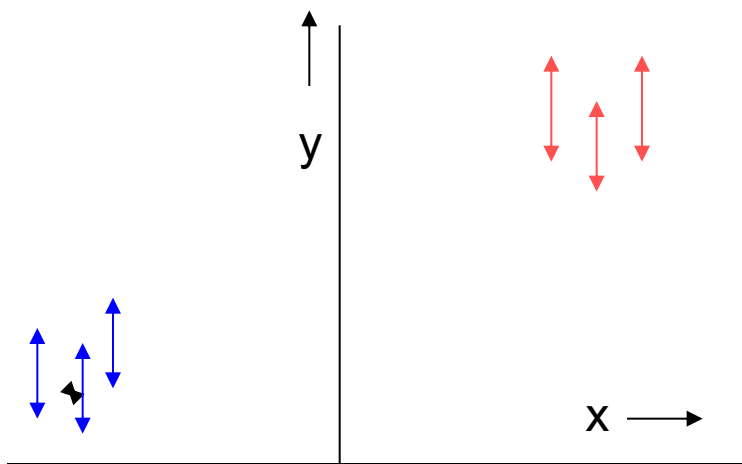COVARIANCE $(a, b) \propto -\langle X \rangle$
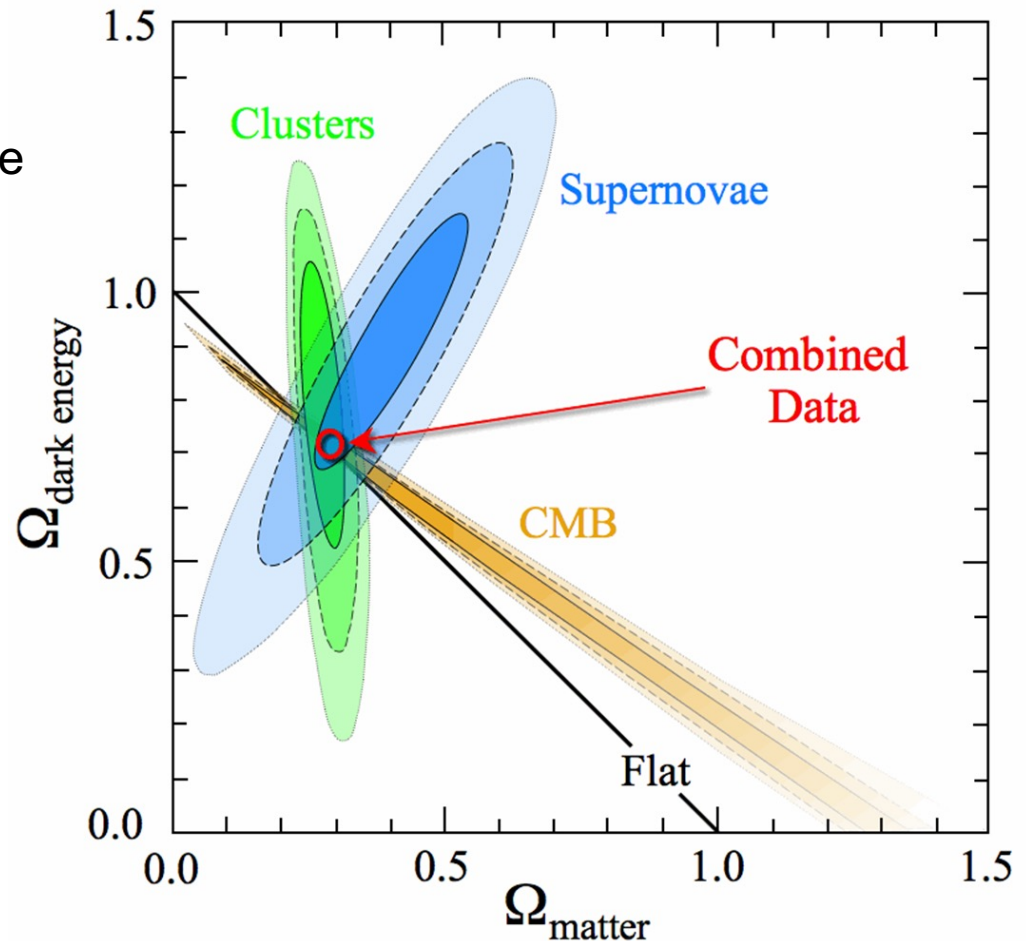
(a)



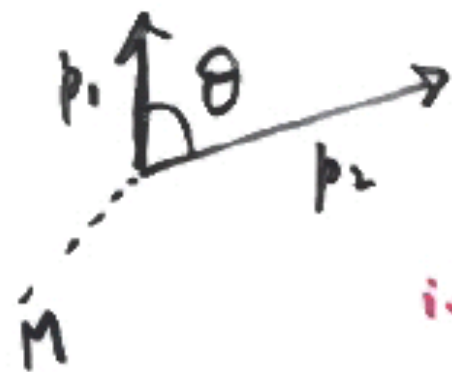$\langle x \rangle$ pos

Fig. 2.4

(b)

$\langle x \rangle$ neg



23

# Uncertainty on $\Omega_{dark\ energy}$

When combining pairs of variables, the uncertainties on the combined parameters can be **much** smaller than any of the individual uncertainties
e.g. $\Omega_{dark\ energy}$

# CORRELATIONS + MASS RESOLUTION

$$M^2 = (E_1 + E_2)^2 - (\underline{p}_1 + \underline{p}_2)^2$$

$$\sim p_1 p_2 \theta \qquad [\, p_i \gg m_i \quad \theta \ll 1 \,]$$

i.e. $M \uparrow$ as $p_i \uparrow$ & $\theta_i \uparrow$

As $p_i \downarrow$, $\theta \uparrow$

$\therefore$ Smaller $\sigma_M$

As $p_i \downarrow$, $\theta \downarrow$

$\therefore$ Larger $\sigma_M$

1) ESTIMATE ERRORS

ESTIMATE CORRELATIONS

(usually easiest if $\rho = 0$ or $\pm 1$)

2) FOR INDEP SOURCES OF ERRORS, ADD ERROR MATRICES

e.g. $M_W$ FROM $WW \rightarrow 4$ JETS

$WW \rightarrow J J l \nu$

$\underline{\underline{E}} = (M_W)_, , (M_W)_2$ ERROR MATRIX

$\underline{\underline{E}} = \underline{\underline{E}}_{stat} + \underline{\underline{E}}_{B.E.} + \underline{\underline{E}}_{E \, scale}$

$+ \underline{\underline{E}}_{FSR} + \underline{\underline{E}}_{colour}$

reconn

$\begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{pmatrix}$

$\begin{pmatrix} \sigma_1^2 & \sigma_1 \sigma_2 \\ \sigma_1 \sigma_2 & \sigma_2^2 \end{pmatrix}$

$\begin{pmatrix} \sigma_1^2 & 0 \\ 0 & 0 \end{pmatrix}$

26

3) TRANSFORMATIONS

e.g. $(x \pm \sigma_x, y \pm \sigma_y)$ with <u>un</u>correl. errors

$\Rightarrow r, \theta$ with correlations

Indep data points

$\Rightarrow$ correlated

a and b

Track fit

4) REPEATED OBSERVATIONS

$(x_i, y_i) \Rightarrow \sigma_x^` \quad \sigma_y^2 \quad$ and

$cov(x, y)$ from $\overline{(x - \bar{x})(y - \bar{y})}$

# Conclusion

Covariance matrix formalism makes life easy when correlations are relevant