

Lecture 4

Multivariate probabilities

Covariance

$$\text{Cov}[X, Y] = E[(X - E(X))(Y - E(Y))] = E[XY] - E[X]E[Y].$$

Correlation

$$\rho = \text{corr}[X, Y] = \frac{\text{Cov}[X, Y]}{\sqrt{V[X]V[Y]}} = \frac{\text{Cov}[X, Y]}{\sigma_X \sigma_Y}.$$

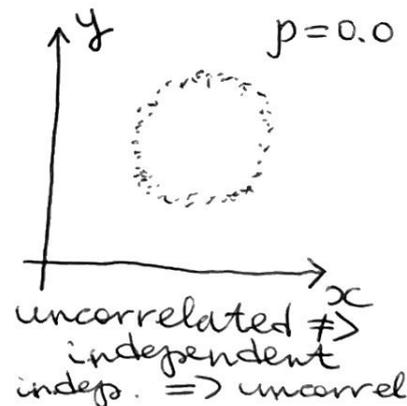
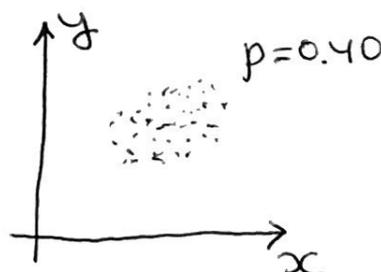
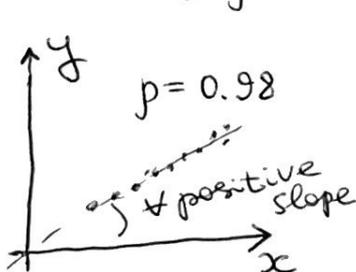
Consider $0 \leq V\left[\frac{X}{\sigma_X} + \frac{Y}{\sigma_Y}\right] = V\left[\frac{X}{\sigma_X}\right] + V\left[\frac{Y}{\sigma_Y}\right] + 2 \text{Cov}\left[\frac{X}{\sigma_X}, \frac{Y}{\sigma_Y}\right] = \frac{V[X]}{\sigma_X^2} + \frac{V[Y]}{\sigma_Y^2} + 2 \frac{\text{Cov}[X, Y]}{\sigma_X \sigma_Y} = 2(1 + \rho),$ or $\rho \geq -1$

likewise, $0 \leq V\left[\frac{X}{\sigma_X} - \frac{Y}{\sigma_Y}\right] = 2(1 - \rho),$ or $\rho \leq 1.$

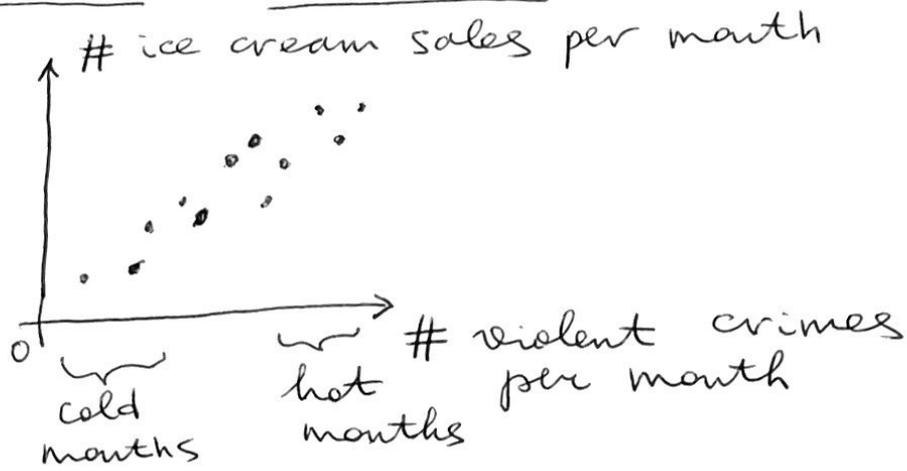
Thus, $\underline{-1 \leq \rho \leq 1}$

Note that $\rho = 1$ iff $Y = aX + b, a > 0$

$[\rho = -1$ iff $Y = aX + b, a < 0]$



Correlation \neq causation

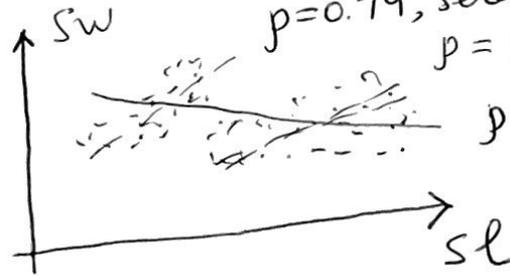


Hidden common cause: weather

storks \sim birth rates
in some areas

Hidden common cause: wealth (?)

Simpson's paradox



$p = 0.74$, setosa

$p = 0.53$ versicolor, $p = 0.46$ virginica

$p = -0.12$, all 3 species

Iris dataset

Multi-D Gaussian distribution

$$\mathcal{N}(\vec{y} | \vec{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(\vec{y}-\vec{\mu})^T \Sigma^{-1}(\vec{y}-\vec{\mu})}$$

$$\vec{y} = (y_1, \dots, y_D)$$

$\Sigma_{D \times D} = \text{Cov}[\vec{y}]$ is the covariance matrix:

$$\Sigma = E[(\vec{y} - \underbrace{\vec{\mu}}_{E[\vec{y}]}) (\vec{y} - \vec{\mu})^T], \text{ or}$$

$$\Sigma = E[\vec{y} \vec{y}^T] - \underline{\underline{\vec{\mu} \cdot \vec{\mu}^T}}.$$

$$\underline{D=2}: \Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12}^2 \\ \sigma_{12}^2 & \sigma_2^2 \end{pmatrix} = \begin{pmatrix} \sigma_1^2 & \rho \sigma_1 \sigma_2 \\ \rho \sigma_1 \sigma_2 & \sigma_2^2 \end{pmatrix}$$

\swarrow corr'n coeff.

Isotropic $\Sigma = \sigma^2 \underbrace{\mathbb{I}_D}_{\text{unit matrix}}$

Marginals and conditionals

Consider $\vec{y} = (\vec{y}_1, \vec{y}_2)$

\uparrow_{2D} \uparrow_D \uparrow_D

assume that ~~the~~ \vec{y} is gaussian distributed:

$$\mathcal{N}(\vec{y} | \vec{\mu}, \Sigma) \Rightarrow \vec{\mu} = \begin{pmatrix} \vec{\mu}_1 \\ \vec{\mu}_2 \end{pmatrix}$$

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}, \quad \Lambda = \Sigma^{-1} = \begin{pmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{pmatrix}$$

\uparrow precision matrix

Marginals: $\int p(\vec{y}_1) = \mathcal{N}(\vec{y}_1 | \vec{\mu}_1, \Sigma_{11}),$
 $\int p(\vec{y}_2) = \mathcal{N}(\vec{y}_2 | \vec{\mu}_2, \Sigma_{22}).$

Moreover,

$$p(\bar{y}_1 | \bar{y}_2) = \mathcal{N}(\bar{y}_1 | \bar{\mu}_{1|2}, \Sigma_{1|2}),$$

where

$$\left\{ \begin{aligned} \bar{\mu}_{1|2} &= \bar{\mu}_1 + \Sigma_{12} \Sigma_{22}^{-1} (\bar{y}_2 - \bar{\mu}_2) = \\ &= \bar{\mu}_1 - \Lambda_{11}^{-1} \Lambda_{12} (\bar{y}_2 - \bar{\mu}_2), \\ \Sigma_{1|2} &= \Lambda_{11}^{-1} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}. \end{aligned} \right.$$

2D example $\sigma^2 \begin{pmatrix} 1 & p \\ p & 1 \end{pmatrix} \Rightarrow \Lambda = \Sigma^{-1} = \begin{pmatrix} 1 & -p \\ -p & 1 \end{pmatrix} \frac{1}{\sigma^2(1-p^2)}$

$$\Sigma = \begin{pmatrix} \sigma^2 & p\sigma^2 \\ p\sigma^2 & \sigma^2 \end{pmatrix} \quad \begin{cases} p(y_1) = \mathcal{N}(y_1 | \mu_1, \sigma^2) \\ p(y_2) = \mathcal{N}(y_2 | \mu_2, \sigma^2) \end{cases}$$

$$p(y_1 | y_2) = ?$$

$$\left\{ \begin{aligned} \mu_{1|2} &= \mu_1 + p\sigma^2 \frac{1}{\sigma^2} (y_2 - \mu_2), \\ \sigma_{1|2} &= \sigma^2 - p\sigma^2 \frac{1}{\sigma^2} p\sigma^2 = \\ &= \sigma^2 [1 - p^2]. \end{aligned} \right.$$

Thus, $p(y_1 | y_2) = \mathcal{N}(y_1 | \mu_{1|2}, \sigma_{1|2}^2) = \mathcal{N}(y_1 | \mu_1 + p(y_2 - \mu_2), \sigma^2(1-p^2))$

Let $\mu_1 = \mu_2 = 0, \sigma = 1, p = 0.8:$

$$E[y_1 | y_2 = 1] = 0.8(1 - 0) = 0.8$$

if y_2 increases by 1.0, from the mean
 y_1 increases by 0.8 on average
 since $p = 0.8$

Next, $V[y_1 | y_2 = 1] = 1 - 0.8^2 = 0.36 < V[y_1] = 1$.

Variance went down b/c observing a correlated variable y_2 provides info about y_1 .

Note that if $\rho = 0$,

$$P(y_1 | y_2) = \mathcal{N}(y_1 | \mu_1, \sigma^2) \stackrel{\text{---}}{=} P(y_1).$$

Observing y_2 does not help since y_1 & y_2 are uncorrelated.

Linear Gaussian systems

$\vec{z} \in \mathbb{R}^L$ hidden variables

$\vec{y} \in \mathbb{R}^D$ noisy observables

Then we can assume:

$$\begin{cases} P(\vec{z}) = \mathcal{N}(\vec{z} | \vec{\mu}_z, \Sigma_z) \\ P(\vec{y} | \vec{z}) = \mathcal{N}(\vec{y} | \mathbf{W}_{D \times L} \vec{z} + \vec{b}_D, \Sigma_y) \end{cases}$$

Thus, the posterior is

$$P(\vec{z} | \vec{y}) = \mathcal{N}(\vec{z} | \vec{\mu}_{z|y}, \Sigma_{z|y}), \text{ where}$$

$$\begin{cases} \Sigma_{z|y}^{-1} = \Sigma_z^{-1} + \mathbf{W}^T \Sigma_y^{-1} \mathbf{W}, \\ \vec{\mu}_{z|y} = \Sigma_{z|y} [\mathbf{W}^T \Sigma_y^{-1} (\vec{y} - \vec{b}) + \Sigma_z^{-1} \vec{\mu}_z] \end{cases}$$

Moreover,

$$p(\bar{y}) = \int d\bar{z} p(\bar{y}|\bar{z}) p(\bar{z}) = \mathcal{N}(\bar{y} | W\bar{\mu}_z + \bar{b}, \Sigma_y + W\Sigma_z W^T).$$

Ex. Bayesian inference of ~~z~~ ^z from

$$\bar{y} = (y_1 \dots y_N) \Rightarrow W = \begin{pmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{pmatrix}_N, \quad \bar{b} = \vec{0}_N$$

$\lambda_0 = \text{precision}$

Assume $p(\bar{z}) = \mathcal{N}(\bar{z} | \mu_0, \lambda_0^{-1})$,
 \uparrow prior $\underbrace{\lambda_0^{-1}}_{\sigma^2}$

$$p(y_i | z) = \mathcal{N}(y_i | z, \lambda^{-1}), \quad i=1, \dots, N$$

$p(\bar{y} | z) = \prod_i p(y_i | z) = \mathcal{N}(\bar{y} | \bar{z}, \lambda^{-1} \mathbb{I}_N)$ \uparrow each y_i is a noisy measurement of z

Then $\mu_{z|y} = \Sigma_{z|y} [N\lambda \bar{y} + \lambda_0 \mu_0] \ominus$

$\Sigma_y^{-1} = \lambda \mathbb{I}_N$ $\uparrow = \lambda \sum_{j=1}^N y_j = N\lambda \bar{y}$
 \uparrow average

$\Sigma_z^{-1} = \lambda_0$

$\Sigma_{z|y}^{-1} = \lambda_0 + N\lambda \equiv \lambda_N$ \leftarrow posterior precision, sum of λ_0 & $N\lambda$

\leftarrow average of y 's

$$\ominus \frac{N\lambda}{N\lambda + \lambda_0} \bar{y} + \frac{\lambda_0}{N\lambda + \lambda_0} \mu_0$$

weighted average of \bar{y} & μ_0

\mathcal{H}_0 $N=1$, $\bar{y} \rightarrow y$ (single datapoint) and

$$\mu_1 = \frac{\lambda_0}{\lambda + \lambda_0} \mu_0 + \frac{\lambda}{\lambda + \lambda_0} y =$$

$\begin{matrix} \text{"} \\ \sigma^{-2} \end{matrix}$
 $\begin{matrix} \text{"} \\ \sigma_0^{-2} \end{matrix}$

$$= \frac{1}{\frac{\sigma_0^2}{\sigma^2} + 1} \mu_0 + \frac{1}{1 + \frac{\sigma^2}{\sigma_0^2}} y = \frac{\sigma^2}{\sigma^2 + \sigma_0^2} \mu_0 + \frac{\sigma_0^2}{\sigma^2 + \sigma_0^2} y$$

$$\textcircled{=} \underbrace{\mu_0}_{\text{prior}} + (y - \mu_0) \underbrace{\frac{\sigma_0^2}{\sigma^2 + \sigma_0^2}}_{\text{correction}} \textcircled{=} \text{due to data}$$

$$\textcircled{=} y - (y - \mu_0) \frac{\sigma^2}{\sigma^2 + \sigma_0^2} \cdot \begin{cases} \sigma_0 \ll \sigma: \mu_1 \approx \mu_0, \\ \text{'prior wins'} \\ \sigma_0 \gg \sigma: \mu_1 \approx y, \\ \text{'data wins'} \end{cases}$$

\uparrow data \uparrow adjustment towards the mean; 'shrinkage'