## 12.2 Covariance Matrices and Joint Probabilities

Linear algebra enters when we run $M$ different experiments at once. We might measure age and height and weight ($M = 3$ measurements of $N$ people). Each experiment has its own mean value. So we have a vector $\boldsymbol{m} = (m_1, m_2, m_3)$ containing the $M$ mean values. Those could be *sample means* of age and height and weight. Or $m_1, m_2, m_3$ could be *expected values* of age, height, weight based on known probabilities.

A matrix becomes involved when we look at variances. Each experiment will have a sample variance $S_i^2$ or an expected $\sigma_i^2 = \text{E}\left[(x_i - m_i)^2\right]$ based on the squared distance from its mean. Those $M$ numbers $\sigma_1^2, \ldots, \sigma_M^2$ will go on the main diagonal of the matrix. So far we have made no connection between the $M$ parallel experiments. They measure $M$ different random variables, but the experiments are not necessarily independent!

If we measure age and height and weight $(\boldsymbol{a}, \boldsymbol{h}, \boldsymbol{w})$ for children, the results will be strongly correlated. Older children are generally taller and heavier. Suppose the means $m_a, m_h, m_w$ are known. Then $\sigma_a^2, \sigma_h^2, \sigma_w^2$ are the separate variances in age, height, weight. **The new numbers are the covariances like $\sigma_{ah}$, where age multiplies height.**

> **Covariance** $\quad \sigma_{ah} = \text{E}\left[(\textbf{age} - \textbf{mean age})\,(\textbf{height} - \textbf{mean height})\right].$ $\quad$ (1)

This definition needs a close look. To compute $\sigma_{ah}$, it is not enough to know the probability of each age and the probability of each height. We have to know the **joint probability of each pair** (**age and height**). This is because age is connected to height.

$\boldsymbol{p_{ah}} =$ probability that a random child has age $= \boldsymbol{a}$ **and** height $= \boldsymbol{h}$: both at once

$\boldsymbol{p_{ij}} = $ **probability that experiment 1 produces $x_i$ and experiment 2 produces $y_j$**

Suppose experiment 1 (age) has mean $m_1$. Experiment 2 (height) has mean $m_2$. The covariance in (1) between experiments 1 and 2 looks at **all pairs** of ages $x_i$, heights $y_j$ :

> **Covariance** $\quad \sigma_{12} = \displaystyle\sum_{\textbf{all }i,\,j} \sum p_{ij}(x_i - m_1)(y_j - m_2)$ $\quad$ (2)

To capture this idea of "joint probability $p_{ij}$" we begin with two small examples.

**Example 1** Flip two coins separately. With 1 for heads and 0 for tails, the results can be $(1, 1)$ or $(1, 0)$ or $(0, 1)$ or $(0, 0)$. Those four outcomes all have probability $p_{11} = p_{10} = p_{01} = p_{00} = \frac{1}{4}$. **Independent experiments have Prob of $(i, j) = $ (Prob of $i$) (Prob of $j$).**

**Example 2** *Glue the coins together*, facing the same way. The only possibilities are $(1, 1)$ and $(0, 0)$. Those have probabilities $\frac{1}{2}$ and $\frac{1}{2}$. The probabilities $p_{10}$ and $p_{01}$ are zero. $(1, 0)$ and $(0, 1)$ won't happen because the coins stick together: both heads or both tails.

**Probability matrices for Examples 1 and 2** $\quad P = \begin{bmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{bmatrix} = \begin{bmatrix} \frac{1}{4} & \frac{1}{4} \\ \frac{1}{4} & \frac{1}{4} \end{bmatrix} \qquad P = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{bmatrix}.$

Let me stay longer with $P$, to show it in good matrix notation. The matrix shows the probability $p_{ij}$ of each pair $(x_i, y_j)$—starting with $(x_1, y_1) =$ (heads, heads) and $(x_1, y_2) =$ (heads, tails). Notice the row sums $p_i$ and column sums $P_j$ and the total sum $= 1$.

$$\textbf{Probability matrix} \quad P = \begin{bmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{bmatrix} \quad \begin{array}{l} p_{11} + p_{12} = \mathbf{p_1} \\ p_{21} + p_{22} = \mathbf{p_2} \end{array} \begin{pmatrix} \text{first} \\ \text{coin} \end{pmatrix}$$

$$\text{(second coin) column sums} \quad \mathbf{P_1} \quad \mathbf{P_2} \qquad \text{4 entries add to 1}$$

Those numbers $p_1, p_2$ and $P_1, P_2$ are called the **marginals** of the matrix $P$:

$p_1 = p_{11} + p_{12} =$ chance of heads from **coin 1** (coin 2 can be heads or tails)
$P_1 = p_{11} + p_{21} =$ chance of heads from **coin 2** (coin 1 can be heads or tails)

Example 1 showed *independent* variables. Every probability $p_{ij}$ equals $p_i$ times $p_j$ $\left(\frac{1}{2} \text{ times } \frac{1}{2} \text{ gave } p_{ij} = \frac{1}{4} \text{ in that example}\right)$. In this case **the covariance $\sigma_{12}$ will be zero**. Heads or tails from the first coin gave no information about the second coin.

$$\boxed{\begin{array}{l} \textbf{Zero covariance } \sigma_{12} \\ \textbf{for independent trials} \end{array} \qquad V = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} = \textbf{diagonal covariance matrix}.}$$

Independent experiments have $\sigma_{12} = 0$ because every $p_{ij} = (p_i)(p_j)$ in equation (2):

$$\sigma_{12} = \sum_i \sum_j (p_i)(p_j)(x_i - m_1)(y_j - m_2) = \left[\sum_i (p_i)(x_i - m_1)\right]\left[\sum_j (p_j)(y_j - m_2)\right] = [\mathbf{0}][\mathbf{0}].$$

The glued coins show perfect correlation. Heads on one means heads on the other. The covariance $\sigma_{12}$ moves from 0 to $\sigma_1 \sigma_2 = \frac{1}{4}$—this is the largest possible value of $\sigma_{12}$:

$$\textbf{Means} = \frac{1}{2} \qquad \sigma_{12} = \frac{1}{2}\left(1 - \frac{1}{2}\right)\left(1 - \frac{1}{2}\right) + 0 + 0 + \frac{1}{2}\left(0 - \frac{1}{2}\right)\left(0 - \frac{1}{2}\right) = \frac{1}{4}$$

Heads or tails from coin 1 gives complete information about heads or tails from coin 2 :

$$\begin{array}{l} \textbf{Glued coins give largest possible covariances} \\ \textbf{Singular covariance matrix: determinant} = 0 \end{array} \qquad V_{\text{glue}} = \begin{bmatrix} \sigma_1^2 & \sigma_1 \sigma_2 \\ \sigma_1 \sigma_2 & \sigma_2^2 \end{bmatrix}$$

**Always $\sigma_1^2 \sigma_2^2 \geq \sigma_{12}^2$**. Thus $\sigma_{12}$ is *between $-\sigma_1\sigma_2$ and $\sigma_1\sigma_2$*. The covariance matrix $V$ is **positive definite** (or in this singular case of glued coins, $V$ is **positive semidefinite**). That is an important fact about $M$ by $M$ covariance matrices for $M$ experiments.

Note that the **sample covariance matrix $S$** from $N$ trials is certainly semidefinite. Every new sample $X = $ (age, height, weight) contributes to the **sample mean** $\overline{X}$ and to $S$. Each term $(X_i - \overline{X})(X_i - \overline{X})^{\mathrm{T}}$ is positive semidefinite and we just add to reach $S$:

$$\overline{X} = \frac{X_1 + \cdots + X_N}{N} \qquad S = \frac{(X_1 - \overline{X})(X_1 - \overline{X})^{\mathrm{T}} + \cdots + (X_N - \overline{X})(X_N - \overline{X})^{\mathrm{T}}}{N - 1} \tag{3}$$

## The Covariance Matrix $V$ is Positive Semidefinite

Come back to the *expected* covariance $\sigma_{12}$ between two experiments 1 and 2 (two coins):

$$
\boxed{
\begin{aligned}
\sigma_{12} &= \text{expected value of } \ [(output\,1 - mean\,1) \text{ times } (output\,2 - mean\,2)] \\
&= \sum_{\text{all } i,\,j} \boldsymbol{p_{ij}} \, (\boldsymbol{x_i} - \boldsymbol{m_1}) \, (\boldsymbol{y_j} - \boldsymbol{m_2}).
\end{aligned}
}
\tag{4}
$$

$p_{ij} \geq 0$ is the probability of seeing output $x_i$ in experiment 1 **and** $y_j$ in experiment 2. Some pair of outputs must appear. Therefore the $N^2$ probabilities $p_{ij}$ add to 1.

$$
\textbf{Total probability (all pairs) is 1} \qquad \sum_{\text{all } i,\,j} p_{ij} = 1.
\tag{5}
$$

Here is another fact we need. *Fix on one particular output $x_i$ in experiment 1. Allow all outputs $y_j$ in experiment 2.* Add the probabilities of $(x_i, y_1), (x_i, y_2), \ldots, (x_i, y_n)$:

$$
\textbf{Row sum } \boldsymbol{p_i} \textbf{ of } \boldsymbol{P} \qquad \sum_{\boldsymbol{j=1}}^{\boldsymbol{n}} \boldsymbol{p_{ij}} = \textbf{probability } \boldsymbol{p_i} \textbf{ of } \boldsymbol{x_i} \textbf{ in experiment 1.}
\tag{6}
$$

Some $y_j$ must happen in experiment 2 ! Whether the two coins are completely separate or glued together, we still get $\frac{1}{2}$ for the probability $p_H = p_{HH} + p_{HT}$ that coin 1 is heads:

$$
\text{(separate) } P_{HH} + P_{HT} = \frac{1}{4} + \frac{1}{4} = \mathbf{\frac{1}{2}} \qquad \text{(glued) } P_{HH} + P_{HT} = \frac{1}{2} + 0 = \mathbf{\frac{1}{2}}.
$$

That basic reasoning allows us to write one matrix formula that includes the covariance $\sigma_{12}$ along with the separate variances $\sigma_1^2$ and $\sigma_2^2$ for experiment 1 and experiment 2. We get the whole covariance matrix $V$ by adding the matrices $V_{ij}$ for each pair $(i, j)$:

$$
\boxed{
\begin{array}{ll}
\textbf{Covariance matrix} \\
\boldsymbol{V = \Sigma\,\Sigma\,V_{ij}}
\end{array}
\quad
V = \sum_{\text{all } i,\,j} p_{ij} \begin{bmatrix} (x_i - m_1)^2 & (x_i - m_1)(y_j - m_2) \\ (x_i - m_1)(y_j - m_2) & (y_j - m_2)^2 \end{bmatrix}
}
\tag{7}
$$

Off the diagonal, this is equation (2) for the covariance $\sigma_{12}$. On the diagonal, we are getting the ordinary variances $\sigma_1^2$ and $\sigma_2^2$. I will show in detail how we get $V_{11} = \sigma_1^2$ by using equation (6). Allowing all $j$ just leaves the probability $p_i$ of $x_i$ in experiment 1 :

$$
\boldsymbol{V_{11}} = \sum_{\text{all } i,\,j} p_{ij}(x_i - m_1)^2 = \sum_{\text{all } i} \text{(probability of } x_i) \, (x_i - m_1)^2 = \boldsymbol{\sigma_1^2}.
\tag{8}
$$

Please look at that twice. It is the key to producing the whole covariance matrix by one formula (7). The beauty of that formula is that it combines 2 by 2 matrices $V_{ij}$. And the matrix $V_{ij}$ in (7) for each pair of outcomes $i, j$ is **positive semidefinite** :

$V_{ij}$ has diagonal entries $p_{ij}(x_i - m_1)^2 \geq 0$ and $p_{ij}(y_j - m_2)^2 \geq 0$ and $\det(V_{ij}) = 0$.

That matrix $V_{ij}$ has rank 1. Equation (7) multiplies $p_{ij}$ *times column $U$ times row $U^{\mathrm{T}}$*:

$$\begin{bmatrix} (x_i - m_1)^2 & (x_i - m_1)(y_j - m_2) \\ (x_i - m_1)(y_j - m_2) & (y_j - m_2)^2 \end{bmatrix} = \begin{bmatrix} x_i - m_1 \\ y_j - m_2 \end{bmatrix} \begin{bmatrix} x_i - m_1 & y_j - m_2 \end{bmatrix} \quad (9)$$

*Every matrix $UU^{\mathrm{T}}$ is positive semidefinite.*  So the whole matrix $V$ (combining these matrices $UU^{\mathrm{T}}$ with weights $p_{ij} \geq 0$) is **at least semidefinite**—and probably $V$ is definite.

**The covariance matrix $V$ is positive definite unless the experiments are dependent.**

Now we move from two variables $x$ and $y$ to $M$ variables like age-height-weight. The output from each trial is a vector $X$ with $M$ components. (Each child has an age-height-weight vector with 3 components.)  The covariance matrix $V$ is now $M$ by $M$. $V$ is created from the output vectors $X$ and their average $\overline{X} = \mathbf{E}\,[X]$ :

**Covariance matrix** $\qquad V = \mathbf{E}\,\left[ \left( X - \overline{X} \right) \left( X - \overline{X} \right)^{\mathrm{T}} \right] \qquad\qquad (10)$

Remember that $XX^{\mathrm{T}}$ and $\overline{X}\,\overline{X}^{\mathrm{T}} = $ (column)(row) are $M$ by $M$ matrices.

For $M = 1$ (one variable) you see that $\overline{X}$ is the mean $m$ and $V$ is $\sigma^2$ (Section 12.1). For $M = 2$ (two coins) you see that $\overline{X}$ is $(m_1, m_2)$ and $V$ matches equation (10).  The expectation E always adds up outputs times their probabilities.  For age-height-weight the output could be $X = $ (5 years, 31 inches, 48 pounds) and its probability is $p_{5,31,48}$ .

Now comes a new idea. *Take any linear combination $c^{\mathrm{T}}X = c_1 X_1 + \cdots + c_M X_M$.* With $c = (6, 2, 5)$ this would be $c^{\mathrm{T}}X = 6\,(\text{age}) + 2\,(\text{height}) + 5\,(\text{weight})$. By linearity we know that its expected value $\mathrm{E}\,[c^{\mathrm{T}}X]$ is $c^{\mathrm{T}}\mathrm{E}\,[X] = c^{\mathrm{T}}\overline{X}$ :

$$\mathbf{E}\,[c^{\mathrm{T}}X] = c^{\mathrm{T}}\mathbf{E}\,[X] = 6\,(\text{expected age}) + 2\,(\text{expected height}) + 5\,(\text{expected weight}).$$

More than that, we also know the *variance $\sigma^2$* of that number $c^{\mathrm{T}}X$:

$$\begin{aligned} \text{Variance of } c^{\mathrm{T}}X &= \mathrm{E}\,\left[ \left( c^{\mathrm{T}}X - c^{\mathrm{T}}\overline{X} \right) \left( c^{\mathrm{T}}X - c^{\mathrm{T}}\overline{X} \right)^{\mathrm{T}} \right] \\ &= c^{\mathrm{T}}\mathrm{E}\,\left[ \left( X - \overline{X} \right) \left( X - \overline{X} \right)^{\mathrm{T}} \right] c = c^{\mathrm{T}}Vc\,! \end{aligned} \qquad (11)$$

Now the key point: *The variance of $c^{\mathrm{T}}X$ can never be negative.* So $c^{\mathrm{T}}Vc \geq 0$. *The covariance matrix $V$ is therefore positive semidefinite by the energy test $c^{\mathrm{T}}Vc \geq 0$.*

Covariance matrices $V$ open up the link between probability and linear algebra: $V$ equals $Q\Lambda Q^{\mathrm{T}}$ with eigenvalues $\lambda_i \geq 0$ and orthonormal eigenvectors $q_1$ to $q_M$.

**Diagonalizing the covariance matrix means finding $M$ *independent* experiments as combinations of the original $M$ experiments.**

**Confession**   I am not entirely happy with that proof based on $c^{\mathrm{T}} V c \geq 0$. The expectation symbol $\mathbf{E}$ is burying the key idea of **joint probability**. Allow me to show directly that $V$ is positive semidefinite (at least for the age-height-weight example). The proof is simply that $V$ **is the sum of the joint probability $p_{ahw}$ of each combination (age, height, weight) times the positive semidefinite matrix $UU^{\mathrm{T}}$**. Here $U$ is $X - \overline{X}$:

$$V = \sum_{\text{all } a,h,w} p_{ahw} \, U \, U^{\mathrm{T}} \quad \text{with} \quad U = \begin{bmatrix} \text{age} \\ \text{height} \\ \text{weight} \end{bmatrix} - \begin{bmatrix} \text{mean age} \\ \text{mean height} \\ \text{mean weight} \end{bmatrix}. \quad (12)$$

This is exactly like the 2 by 2 coin flip matrix $V$ in equation (7). Now $M = 3$.

The value of the expectation symbol E is that it also allows *pdf*'s (probability density functions like $p(x, y, z)$ for continuous random variables $x$ and $y$ and $z$). If we allow all numbers as ages and heights and weights, instead of age $i = 0, 1, 2, 3 \ldots$, then we need $p(x, y, z)$ instead of $p_{ijk}$. The sums in this section of the book would all change to integrals. But we still have $V = \mathrm{E}\left[ UU^{\mathrm{T}} \right]$:

**Covariance matrix** $\quad V = \displaystyle\iiint p(x, y, z) \, UU^{\mathrm{T}} \, dx \, dy \, dz \quad \text{with} \quad U = \begin{bmatrix} x - \overline{x} \\ y - \overline{y} \\ z - \overline{z} \end{bmatrix}. \quad (13)$

Always $\iiint p = 1$. Examples 1–2 emphasized how $p$ can give diagonal $V$ or singular $V$:

   **Independent variables $x, y, z$** $\quad p(x, y, z) = p_1(x) \, p_2(y) \, p_3(z)$.

   **Dependent  variables  $x, y, z$** $\quad p(x, y, z) = 0$ except when $cx + dy + ez = 0$.

## The Mean and Variance of $z = x + y$

Start with the sample mean. We have $N$ samples of $x$. Their mean (= average) is $m_x$. We also have $N$ samples of $y$ and their mean is $m_y$. **The sample mean of $z = x + y$ is clearly $m_z = m_x + m_y$**:

**Mean of sum = Sum of means** $\quad \dfrac{1}{N} \displaystyle\sum_1^N (x_i + y_i) = \dfrac{1}{N} \sum_1^N x_i + \dfrac{1}{N} \sum_1^N y_i. \quad (14)$

Nice to see something that simple. The *expected* mean of $z = x + y$ doesn't look so simple, but it must come out as $\mathbf{E}[z] = \mathbf{E}[x] + \mathbf{E}[y]$. Here is one way to see this.

The joint probability of the pair $(x_i, y_j)$ is $p_{ij}$. Its value depends on whether the experiments are independent, which we don't know. But for the mean of the sum $z = x + y$,

dependence or independence of $x$ and $y$ doesn't matter. *Expected values still add*:

$$\mathbf{E}[\boldsymbol{x} + \boldsymbol{y}] = \sum_i \sum_j p_{ij}(x_i + y_j) = \sum_i \sum_j p_{ij}x_i + \sum_i \sum_j p_{ij}y_j. \qquad (15)$$

All the sums go from $1$ to $N$. We can add in any order. For the first term on the right side, add the $p_{ij}$ along row $i$ of the probability matrix $P$ to get $p_i$. That double sum gives $\mathrm{E}[x]$:

$$\sum_i \sum_j p_{ij}x_i = \sum_i (p_{i1} + \cdots + p_{iN})x_i = \sum_i p_i x_i = \mathrm{E}[x].$$

For the last term, add $p_{ij}$ down column $j$ of the matrix to get the probability $P_j$ of $y_j$. Those pairs $(x_1, y_j)$ and $(x_2, y_j)$ and $\ldots$ and $(x_N, y_j)$ are all the ways to produce $y_j$:

$$\sum_i \sum_j p_{ij}y_j = \sum_j (p_{1j} + \cdots + p_{Nj})y_j = \sum_j P_j y_j = \mathrm{E}[y].$$

Now equation (15) says that $\mathrm{E}[\boldsymbol{x} + \boldsymbol{y}] = \mathrm{E}[\boldsymbol{x}] + \mathrm{E}[\boldsymbol{y}]$.

What about the variance of $z = x + y$? The joint probabilities $p_{ij}$ and the covariance $\sigma_{xy}$ will be involved. Let me separate the variance of $x + y$ into three simple pieces:

$$\boldsymbol{\sigma_z^2} = \sum\sum p_{ij}(x_i + y_j - m_x - m_y)^2$$
$$= \sum\sum p_{ij}(x_i - m_x)^2 + \sum\sum p_{ij}(y_j - m_y)^2 + 2\sum\sum p_{ij}(x_i - m_x)(y_j - m_y)$$

The first piece is $\boldsymbol{\sigma_x^2}$. The second piece is $\boldsymbol{\sigma_y^2}$. The last piece is $\boldsymbol{2\sigma_{xy}}$.

$$\textbf{The variance of } z = x + y \quad \textbf{is} \quad \sigma_z^2 = \sigma_x^2 + \sigma_y^2 + 2\sigma_{xy}. \qquad (16)$$

## The Covariance Matrix for $Z = AX$

Here is a good way to see $\sigma_z^2$ when $z = x + y$. Think of $(x, y)$ as a column vector $\boldsymbol{X}$. Think of the 1 by 2 matrix $A = \begin{bmatrix} 1 & 1 \end{bmatrix}$ multiplying that vector $\boldsymbol{X}$. Then $A\boldsymbol{X}$ is the sum $z = x + y$. The variance $\boldsymbol{\sigma_z^2}$ in equation (16) goes into matrix notation as

$$\boldsymbol{\sigma_z^2} = \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_y^2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \text{which is} \quad \boldsymbol{\sigma_z^2 = AVA^{\mathrm{T}}}. \qquad (17)$$

You can see that $\sigma_z^2 = AVA^{\mathrm{T}}$ in (17) agrees with $\sigma_x^2 + \sigma_y^2 + 2\sigma_{xy}$ in (16).

Now for the main point. The vector $\boldsymbol{X}$ could have $M$ components coming from $M$ experiments (instead of only 2). Those experiments will have an $M$ by $M$ covariance matrix $\boldsymbol{V_X}$. The matrix $A$ could be $K$ by $M$. Then $A\boldsymbol{X}$ is a vector with $K$ combinations of the $M$ outputs (instead of 1 combination $x + y$ of 2 outputs).

That vector $\boldsymbol{Z} = A\boldsymbol{X}$ of length $K$ has a $K$ by $K$ covariance matrix $V_{\boldsymbol{Z}}$. Then the great rule for covariance matrices—of which equation (17) was only a 1 by 2 example—is this beautiful formula: Covariance matrix of $A\boldsymbol{X}$ is $A$ (covariance matrix of $\boldsymbol{X}$) $A^{\mathrm{T}}$:

$$\boxed{\textbf{The covariance matrix of } \boldsymbol{Z} = A\boldsymbol{X} \textbf{ is } \quad V_{\boldsymbol{Z}} = AV_{\boldsymbol{X}}A^{\mathrm{T}}} \qquad (18)$$

To me, this neat formula shows the beauty of matrix multiplication. I won't prove this formula, just admire it. It is constantly used in applications—coming in Section 12.3.

### The Correlation $\rho$

Correlation $\rho_{xy}$ is closely related to covariance $\sigma_{xy}$. They both measure dependence or independence. Start by rescaling or "standardizing" the random variables $x$ and $y$ **The new $X = x/\sigma_x$ and $Y = y/\sigma_y$ have variance $\sigma_X^2 = \sigma_Y^2 = 1$**. This is just like dividing a vector $v$ by its length to produce a unit vector $v/||v||$ of length 1.

**The correlation of $x$ and $y$ is the covariance of $X$ and $Y$**. If the original covariance of $x$ and $y$ was $\sigma_{xy}$, then rescaling to $X$ and $Y$ will divide by $\sigma_x$ and $\sigma_y$ :

$$\text{Correlation } \rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} = \text{covariance of } \frac{x}{\sigma_x} \text{ and } \frac{y}{\sigma_y} \qquad \text{Always } -1 \le \rho_{xy} \le 1$$

Zero covariance gives zero correlation. *Independent random variables* produce $\rho_{xy} = 0$.

We know that always $\sigma_{xy}^2 \le \sigma_x^2 \sigma_y^2$ (the covariance matrix $V$ is at least positive semidefinite). Then $\rho_{xy}^2 \le 1$. Correlation near $\rho = +1$ means strong dependence in the same direction : often voting the same. Negative correlation means that $y$ tends to be below its mean when $x$ is above its mean : Voting in opposite directions.

**Example 3** *Suppose that $y$ is just $-x$.* A coin flip has outputs $x = 0$ or 1. The same flip has outputs $y = 0$ or $-1$. The mean $m_x$ is $\frac{1}{2}$ for a fair coin, and $m_y$ is $-\frac{1}{2}$. The covariance is $\sigma_{xy} = -\sigma_x \sigma_y$. The correlation divides by $\sigma_x \sigma_y$ to get $\rho_{xy} = -1$. In this case the correlation matrix $R$ has determinant zero (singular and only semidefinite) :

$$\text{Correlation matrix } \quad R = \begin{bmatrix} 1 & \rho_{xy} \\ \rho_{xy} & 1 \end{bmatrix} \qquad R = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \text{ when } y = -x$$

$R$ always has 1's on the diagonal because we normalized to $\sigma_X = \sigma_Y = 1$. $R$ is the correlation matrix for $x$ and $y$, and the covariance matrix for $X = x/\sigma_x$ and $Y = y/\sigma_y$.

That number $\rho_{xy}$ is also called the Pearson coefficient.

**Example 4** Suppose the random variables $x, y, z$ are *independent. What matrix is $R$?*

*Answer* $R$ *is the identity matrix.* All three correlations $\rho_{xx}, \rho_{yy}, \rho_{zz}$ are 1 by definition. All three cross-correlations $\rho_{xy}, \rho_{xz}, \rho_{yz}$ are zero by independence.

The correlation matrix $R$ comes from the covariance matrix $V$, when we rescale every row and every column. Divide each row $i$ and column $i$ by the $i$th standard deviation $\sigma_i$.

(a) $R = DVD$ for the diagonal matrix $D = \text{diag } [1/\sigma_1, \ldots, 1/\sigma_M]$.

(b) If covariance $V$ is positive definite, correlation $R = DVD$ is also positive definite.

■   **WORKED EXAMPLES**   ■

**12.2 A**    Suppose $x$ and $y$ are independent random variables with mean $0$ and variance $1$. Then the covariance matrix $V_X$ for $X = (x, y)$ is the 2 by 2 identity matrix. What are the mean $m_Z$ and the covariance matrix $V_Z$ for the 3-component vector $Z = (x, y, ax + by)$ ?

**Solution**

$$\textbf{Z is connected to X by A} \quad Z = \begin{bmatrix} x \\ y \\ ax + by \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ a & b \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = AX.$$

The vector $m_X$ contains the means of the $M$ components of $X$. The vector $m_Z$ contains the means of the $K$ components of $Z = AX$. The matrix connection between the means of $X$ and $Z$ has to be linear: $m_Z = A\,m_X$. The mean of $ax + by$ is $am_x + bm_y$.

     The covariance matrix for $Z$ is $V_Z = AA^T$, when $V_X$ is the 2 by 2 identity matrix:

$$V_Z = \begin{array}{c} \textbf{covariance matrix for} \\ \boldsymbol{Z = (x, y, ax + by)} \end{array} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ a & b \end{bmatrix} \begin{bmatrix} 1 & 0 & a \\ 0 & 1 & b \end{bmatrix} = \begin{bmatrix} \mathbf{1} & \mathbf{0} & \mathbf{a} \\ \mathbf{0} & \mathbf{1} & \mathbf{b} \\ \mathbf{a} & \mathbf{b} & \mathbf{a^2 + b^2} \end{bmatrix}.$$

Interpretation: $x$ and $y$ are independent so $\sigma_{xy} = 0$. Then the covariance of $x$ with $ax + by$ is $a$ and the covariance of $y$ with $ax + by$ is $b$. Those just come from the two independent parts of $ax + by$. Finally, equation (18) gives the variance of $ax + by$:

$$\textbf{Use } V_Z = A V_X A^T \qquad \sigma^2_{ax+by} = \sigma^2_{ax} + \sigma^2_{by} + 2\sigma_{ax,by} = a^2 + b^2 + 0.$$

The 3 by 3 matrix $V_Z$ is *singular*. Its determinant is $a^2 + b^2 - a^2 - b^2 = 0$. The third component $z = ax + by$ is completely dependent on $x$ and $y$. The rank of $V_Z$ is only 2.

**GPS Example**    The signal from a GPS satellite includes its departure time. The receiver clock gives the arrival time. The receiver multiplies the travel time by the speed of light. Then it knows the distance from that satellite. Distances from four or more satellites pinpoint the receiver position (using least squares !).

     One problem: The speed of light changes in the ionosphere. But the correction will be almost the same for all nearby receivers. If one receiver stays in a known position, we can take differences from that position. **Differential GPS** reduces the error variance:

$$\begin{array}{cc} \textbf{Difference matrix} & \textbf{Covariance matrix} \\ A = \begin{bmatrix} 1 & -1 \end{bmatrix} & V_Z = A V_X A^T \end{array} \quad \begin{aligned} V_Z &= \begin{bmatrix} 1 & -1 \end{bmatrix} \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} \\ &= \sigma_1^2 - 2\sigma_{12} + \sigma_2^2 \end{aligned}$$

Errors in the speed of light are gone. Then centimeter positioning accuracy is achievable. (The key ideas are on page 320 of *Algorithms for Global Positioning* by Borre and Strang.) The GPS world is all about time and space and amazing accuracy.

## Problem Set 12.2

**1**    (a) Compute the variance $\sigma^2$ when the coin flip probabilities are $p$ and $1 - p$ (tails $= 0$, heads $= 1$).

       (b) The sum of $N$ independent flips (0 or 1) is the count of heads after $N$ tries. The rule (16-17-18) for the variance of a sum gives $\sigma^2 = $ _____ .

**2**    What is the covariance $\sigma_{kl}$ between the results $x_1, \ldots, x_n$ of Experiment 3 and the results $y_1, \ldots, y_n$ of Experiment 5 ? Your formula will look like $\sigma_{12}$ in equation (2). Then the $(3, 5)$ and $(5, 3)$ entries of the covariance matrix $V$ are $\sigma_{35} = \sigma_{53}$.

**3**    For $M = 3$ experiments, the variance-covariance matrix $V$ will be 3 by 3. There will be a probability $p_{ijk}$ that the three outputs are $x_i$ and $y_j$ and $z_k$. Write down a formula like equation (7) for the matrix $V$.

**4**    What is the covariance matrix $V$ for $M = 3$ independent experiments with means $m_1, m_2, m_3$ and variances $\sigma_1^2, \sigma_2^2, \sigma_3^2$ ?

**Problems 5–9 are about the conditional probability that $Y = y_j$ when we know $X = x_i$.**
Notation: **Prob** $(Y = y_j | X = x_i) = $ probability of the outcome $y_j$ given that $X = x_i$.

*Example* 1    *Coin* 1 *is glued to coin* 2. Then  Prob $(Y = $ heads when $X = $ heads) is **1**.
*Example* 2    *Independent coin flips* : $X$ gives no information about $Y$. Useless to know $X$.
       Then Prob $(Y = $ heads $| X = $ heads) is the same as Prob $(Y = $ heads).

**5**    Explain the **sum rule** of conditional probability :

$$\text{Prob} \, (Y = y_j) = \text{ sum over all outputs } x_i \text{ of Prob} \, (Y = y_j | X = x_i).$$

**6**    The $n$ by $n$ matrix $P$ contains **joint probabilities** $p_{ij} = $ Prob $(X = x_i$ **and** $Y = y_j)$.

       Explain why the conditional Prob $(Y = y_j | X = x_i)$ equals $\dfrac{p_{ij}}{p_{i1} + \cdots + p_{in}} = \dfrac{p_{ij}}{p_i}$.

**7**    For this joint probability matrix with Prob $(x_1, y_2) = 0.3$, find Prob $(y_2 | x_1)$ and Prob $(x_1)$.

$$P = \begin{bmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{bmatrix} = \begin{bmatrix} 0.1 & 0.3 \\ 0.2 & 0.4 \end{bmatrix} \qquad \begin{array}{l} \text{The entries } p_{ij} \text{ add to 1.} \\ \text{Some } i, j \text{ must happen.} \end{array}$$

**8**    Explain the **product rule** of conditional probability:

$$p_{ij} = \text{Prob} \, (X = x_i \text{ **and** } Y = y_j) \text{ equals Prob} \, (Y = y_j | X = x_i) \text{ times Prob} \, (X = x_i).$$

**9**    Derive this **Bayes Theorem** for $p_{ij}$ from the product rule in Problem 8:

$$\text{Prob} \, (Y = y_j \text{ **and** } X = x_i) = \frac{\text{Prob} \, (X = x_i | Y = y_j) \, \text{Prob} \, (Y = y_j)}{\text{Prob} \, (X = x_i)}$$

"Bayesians" use prior information. "Frequentists" only use sampling information.