

# Chapter 12

## Linear Algebra in Probability & Statistics

### 12.1 Mean, Variance, and Probability

We are starting with the three fundamental words of this chapter: *mean*, *variance*, and *probability*. Let me give a rough explanation of their meaning before I write any formulas:

The **mean** is the *average value* or expected value

The **variance**  $\sigma^2$  measures the average *squared distance* from the mean  $m$

The **probabilities** of  $n$  different outcomes are positive numbers  $p_1, \dots, p_n$  adding to 1.

Certainly the mean is easy to understand. We will start there. But right away we have two different situations that you have to keep straight. On the one hand, we may have the results (*sample values*) from a completed trial. On the other hand, we may have the expected results (*expected values*) from future trials. Let me give examples:

**Sample values** Five random freshmen have ages **18, 17, 18, 19, 17**

**Sample mean**  $\frac{1}{5}(18 + 17 + 18 + 19 + 17) = \mathbf{17.8}$

**Probabilities** The ages in a freshmen class are 17 (**20%**), 18 (**50%**), 19 (**30%**)

A random freshman has **expected age**  $E[x] = (0.2)17 + (0.5)18 + (0.3)19 = \mathbf{18.1}$

Both numbers 17.8 and 18.1 are correct averages. The sample mean starts with  $N$  samples  $x_1, \dots, x_N$  from a completed trial. Their mean is the *average* of the  $N$  observed samples:

$$\text{Sample mean} \quad m = \mu = \frac{1}{N}(x_1 + x_2 + \dots + x_N) \quad (1)$$

The **expected value of  $x$**  starts with the probabilities  $p_1, \dots, p_n$  of the ages  $x_1, \dots, x_n$ :

$$\text{Expected value } m = E[x] = p_1x_1 + p_2x_2 + \dots + p_nx_n. \quad (2)$$

This is  $p \cdot x$ . Notice that  $m = E[x]$  tells us what to expect,  $m = \mu$  tells us what we got.

By taking many samples (large  $N$ ), the sample results will come close to the probabilities. The “Law of Large Numbers” says that with probability 1, the sample mean will converge to its expected value  $E[x]$  as the sample size  $N$  increases. A fair coin has probability  $p_0 = \frac{1}{2}$  of tails and  $p_1 = \frac{1}{2}$  of heads. Then  $E[x] = (\frac{1}{2})0 + \frac{1}{2}(1)$ . The fraction of heads in  $N$  flips of the coin is the sample mean, expected to approach  $E[x] = \frac{1}{2}$ .

This does *not* mean that if we have seen more tails than heads, the next sample is likely to be heads. The odds remain 50-50. The first 100 or 1000 flips do affect the sample mean. *But 1000 flips will not affect its limit*—because you are dividing by  $N \rightarrow \infty$ .

### Variance (around the mean)

The **variance  $\sigma^2$**  measures expected distance (squared) from the expected mean  $E[x]$ . The **sample variance  $S^2$**  measures actual distance (squared) from the sample mean. The square root is the **standard deviation  $\sigma$  or  $S$** . After an exam, I email  $\mu$  and  $S$  to the class. I don’t know the expected mean and variance because I don’t know the probabilities  $p_1$  to  $p_{100}$  for each score. (After teaching for 50 years, I still have no idea what to expect.)

The deviation is always deviation *from the mean*—sample or expected. We are looking for the size of the “spread” around the mean value  $x = m$ . Start with  $N$  samples.

$$\text{Sample variance } S^2 = \frac{1}{N-1} [(x_1 - m)^2 + \dots + (x_N - m)^2] \quad (3)$$

The sample ages  $x = 18, 17, 18, 19, 17$  have mean  $m = 17.8$ . That sample has variance 0.7:

$$S^2 = \frac{1}{4} [(.2)^2 + (-.8)^2 + (.2)^2 + (1.2)^2 + (-.8)^2] = \frac{1}{4}(2.8) = 0.7$$

The minus signs disappear when we compute squares. Please notice! Statisticians divide by  $N - 1 = 4$  (and not  $N = 5$ ) so that  $S^2$  is an unbiased estimate of  $\sigma^2$ . One degree of freedom is already accounted for in the sample mean.

An important identity comes from splitting each  $(x - m)^2$  into  $x^2 - 2mx + m^2$ :

$$\begin{aligned} \text{sum of } (x_i - m)^2 &= (\text{sum of } x_i^2) - 2m(\text{sum of } x_i) + (\text{sum of } m^2) \\ &= (\text{sum of } x_i^2) - 2m(Nm) + Nm^2 \\ \text{sum of } (x_i - m)^2 &= (\text{sum of } x_i^2) - Nm^2. \end{aligned} \quad (4)$$

This is an equivalent way to find  $(x_1 - m)^2 + \dots + (x_N - m)^2$  by adding  $x_1^2 + \dots + x_N^2$ .

Now start with probabilities  $p_i$  (never negative !) instead of samples. We find expected values instead of sample values. The variance  $\sigma^2$  is the crucial number in statistics.

$$\text{Variance } \sigma^2 = E[(x - m)^2] = p_1(x_1 - m)^2 + \cdots + p_n(x_n - m)^2. \quad (5)$$

We are squaring the distance from the expected value  $m = E[x]$ . We don't have samples, only expectations. We know probabilities but we don't know experimental outcomes.

**Example 1** Find the variance  $\sigma^2$  of the ages of college freshmen.

**Solution** The probabilities of ages  $x_i = 17, 18, 19$  were  $p_i = 0.2$  and  $0.5$  and  $0.3$ . The expected value was  $m = \sum p_i x_i = 18.1$ . The variance uses those same probabilities:

$$\begin{aligned} \sigma^2 &= (0.2)(17 - 18.1)^2 + (0.5)(18 - 18.1)^2 + (0.3)(19 - 18.1)^2 \\ &= (0.2)(1.21) + (0.5)(0.01) + (0.3)(0.81) = 0.49. \end{aligned}$$

The **standard deviation** is the square root  $\sigma = 0.7$ .

This measures the spread of 17, 18, 19 around  $E[x]$ , weighted by probabilities .2, .5, .3.

## Continuous Probability Distributions

Up to now we have allowed for  $n$  possible outcomes  $x_1, \dots, x_n$ . With ages 17, 18, 19, we only had  $n = 3$ . If we measure age in days instead of years, there will be a thousand possible ages (too many). Better to allow *every number between 17 and 20*—a continuum of possible ages. Then the probabilities  $p_1, p_2, p_3$  for ages  $x_1, x_2, x_3$  have to move to a **probability distribution**  $p(x)$  for a whole continuous range of ages  $17 \leq x \leq 20$ .

The best way to explain probability distributions is to give you two examples. They will be the **uniform distribution** and the **normal distribution**. The first (uniform) is easy. The normal distribution is all-important.

**Uniform distribution** Suppose ages are uniformly distributed between 17.0 and 20.0. All ages between those numbers are “equally likely”. Of course any one exact age has no chance at all. There is zero probability that you will hit the exact number  $x = 17.1$  or  $x = 17 + \sqrt{2}$ . What you can truthfully provide (assuming our uniform distribution) is **the chance  $F(x)$  that a random freshman has age less than  $x$** :

The chance of age less than  $x = 17$  is  $F(17) = 0$        $x \leq 17$  won't happen  
 The chance of age less than  $x = 20$  is  $F(20) = 1$        $x \leq 20$  will happen  
 The chance of age less than  $x$  is  $F(x) = \frac{1}{3}(x - 17)$        $F$  goes from 0 to 1

That formula  $F(x) = \frac{1}{3}(x - 17)$  gives  $F = 0$  at  $x = 17$ ; then  $x \leq 17$  won't happen. It gives  $F(x) = 1$  at  $x = 20$ ; then  $x \leq 20$  is sure. Between 17 and 20, the graph of the **cumulative distribution**  $F(x)$  increases linearly for this uniform model.

Let me draw the graphs of  $F(x)$  and its derivative  $p(x)$  = “probability density function”.

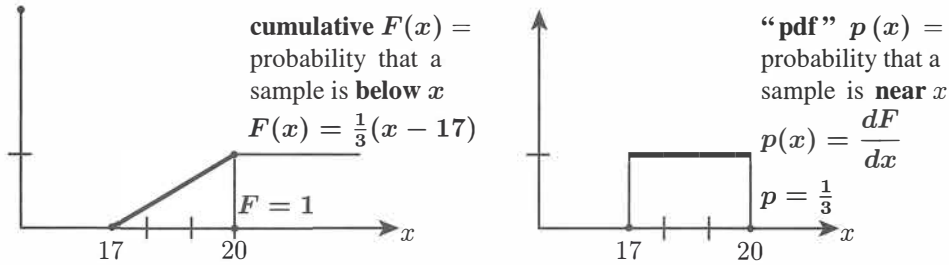


Figure 12.1:  $F(x)$  is the cumulative distribution and its derivative  $p(x) = dF/dx$  is the **probability density function (pdf)**. For this uniform distribution,  $p(x)$  is constant between 17 and 20. The total area under the graph of  $p(x)$  is the total probability  $F = 1$ .

You could say that  $p(x) dx$  is the probability of a sample falling in between  $x$  and  $x + dx$ . This is “infinitesimally true”:  $p(x) dx$  is  $F(x + dx) - F(x)$ . Here is the full truth:

$$F = \text{integral of } p \quad \text{Probability of } a \leq x \leq b = \int_a^b p(x) dx = F(b) - F(a) \quad (6)$$

$F(b)$  is the probability of  $x \leq b$ . I subtract  $F(a)$  to keep  $x \geq a$ . That leaves  $a \leq x \leq b$ .

### Mean and Variance of $p(x)$

What are the mean  $m$  and variance  $\sigma^2$  for a probability distribution? Previously we added  $p_i x_i$  to get the mean (expected value). With a continuous distribution we **integrate**  $x p(x)$ :

$$\text{Mean} \quad m = E[x] = \int x p(x) dx = \int_{x=17}^{20} (x) \left(\frac{1}{3}\right) dx = 18.5$$

For this uniform distribution, the mean  $m$  is halfway between 17 and 20. Then the probability of a random value  $x$  below this halfway point  $m = 18.5$  is  $F(m) = \frac{1}{2}$ .

In MATLAB,  $x = \text{rand}(1)$  chooses a random number uniformly between 0 and 1. Then the expected mean is  $m = \frac{1}{2}$ . The interval from 0 to  $x$  has probability  $F(x) = x$ . The interval below the mean  $m$  always has probability  $F(m) = \frac{1}{2}$ .

The variance is the average squared distance to the mean. With  $N$  outcomes,  $\sigma^2$  is the sum of  $p_i (x_i - m)^2$ . For a continuous random variable  $x$ , the sum changes to an **integral**.

$$\text{Variance} \quad \sigma^2 = E[(x - m)^2] = \int p(x) (x - m)^2 dx \quad (7)$$

When ages are uniform between  $17 \leq x \leq 20$ , the integral can shift to  $0 \leq x \leq 3$ :

$$\sigma^2 = \int_{17}^{20} \frac{1}{3}(x - 18.5)^2 dx = \int_0^3 \frac{1}{3}(x - 1.5)^2 dx = \frac{1}{9}(x - 1.5)^3 \bigg|_{x=0}^{x=3} = \frac{2}{9}(1.5)^3 = \frac{3}{4}.$$

That is a typical example, and here is the complete picture for a uniform  $p(x)$ , 0 to  $a$ .

**Uniform distribution for  $0 \leq x \leq a$**     **Density**  $p(x) = \frac{1}{a}$     **Cumulative**  $F(x) = \frac{x}{a}$

**Mean**  $m = \frac{a}{2}$  halfway    **Variance**  $\sigma^2 = \int_0^a \frac{1}{a} \left(x - \frac{a}{2}\right)^2 dx = \frac{a^2}{12}$     (8)

The mean is a multiple of  $a$ , the variance is a multiple of  $a^2$ . For  $a = 3$ ,  $\sigma^2 = \frac{9}{12} = \frac{3}{4}$ . For one random number between 0 and 1 (mean  $\frac{1}{2}$ ) the variance is  $\sigma^2 = \frac{1}{12}$ .

### Normal Distribution : Bell-shaped Curve

The normal distribution is also called the “Gaussian” distribution. It is the most important of all probability density functions  $p(x)$ . The reason for its overwhelming importance comes from repeating an experiment and averaging the outcomes. The experiments have their own distribution (like heads and tails). *The average approaches a normal distribution.*

**Central Limit Theorem (informal)** The average of  $N$  samples of “any” probability distribution approaches a normal distribution as  $N \rightarrow \infty$ .

Start with the “standard normal distribution”. It is symmetric around  $x = 0$ , so its mean value is  $m = 0$ . It is chosen to have a standard variance  $\sigma^2 = 1$ . It is called  $N(0, 1)$ .

|  |     |
|--|-----|
| <b>Standard normal distribution</b> $p(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}.$ | (9) |
|--|-----|

The graph of  $p(x)$  is the **bell-shaped curve** in Figure 12.2. The standard facts are

|                              |   |
|------------------------------|---|
| <b>Total probability</b> = 1 | $\int_{-\infty}^{\infty} p(x) dx = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-x^2/2} dx = 1$ |
| <b>Mean</b> $E[x] = 0$       | $m = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x e^{-x^2/2} dx = 0$                             |
| <b>Variance</b> $E[x^2] = 1$ | $\sigma^2 = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} (x - 0)^2 e^{-x^2/2} dx = 1$              |

The zero mean was easy because we are integrating an odd function. Changing  $x$  to  $-x$  shows that “integral =  $-$  integral”. So that integral must be  $m = 0$ .

The other two integrals apply the idea in Problem 12 to reach 1. Figure 12.2 shows a graph of  $p(x)$  for the normal distribution  $\mathbf{N}(0, \sigma)$  and also its cumulative distribution  $F(x) = \text{integral of } p(x)$ . From the symmetry of  $p(x)$  you see *mean = zero*. From  $F(x)$  you see a very important practical approximation for opinion polling:

The probability that a random sample falls between  $-\sigma$  and  $\sigma$  is  $F(\sigma) - F(-\sigma) \approx \frac{2}{3}$ .

This is because  $\int_{-\sigma}^{\sigma} p(x) dx$  equals  $\int_{-\infty}^{\sigma} p(x) dx - \int_{-\infty}^{-\sigma} p(x) dx = F(\sigma) - F(-\sigma)$ .

Similarly, the probability that a random  $x$  lies between  $-2\sigma$  and  $2\sigma$  (“*less than two standard deviations from the mean*”) is  $F(2\sigma) - F(-2\sigma) \approx 0.95$ . If you have an experimental result further than  $2\sigma$  from the mean, it is fairly sure to be not accidental: chance = 0.05. Drug tests may look for a tighter confirmation, like probability 0.001. Searching for the Higgs boson used a hyper-strict test of  $5\sigma$  deviation from pure accident.

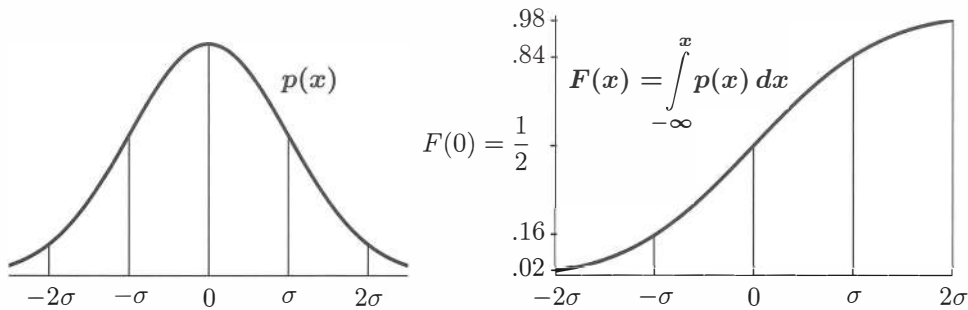


Figure 12.2: The standard normal distribution  $p(x)$  has mean  $m = 0$  and  $\sigma = 1$ .

The normal distribution with any mean  $m$  and standard deviation  $\sigma$  comes by shifting and stretching the standard  $\mathbf{N}(0, 1)$ . **Shift  $x$  to  $x - m$ . Stretch  $x - m$  to  $(x - m)/\sigma$ .**

**Gaussian density  $p(x)$**

**Normal distribution  $\mathbf{N}(m, \sigma)$**

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-m)^2/2\sigma^2} \quad (10)$$

The integral of  $p(x)$  is  $F(x)$ —the probability that a random sample will fall below  $x$ . The differential  $p(x) dx = F(x + dx) - F(x)$  is the probability that a random sample will fall between  $x$  and  $x + dx$ . There is no simple formula to integrate  $e^{-x^2/2}$ , so this cumulative distribution  $F(x)$  is computed and tabulated very carefully.

**$N$  Coin Flips and  $N \rightarrow \infty$** 

**Example 2** Suppose  $x$  is 1 or  $-1$  with equal probabilities  $p_1 = p_{-1} = \frac{1}{2}$ .

The mean value is  $m = \frac{1}{2}(1) + \frac{1}{2}(-1) = 0$ . The variance is  $\sigma^2 = \frac{1}{2}(1)^2 + \frac{1}{2}(-1)^2 = 1$ .

The key question is the *average*  $A_N = (x_1 + \cdots + x_N)/N$ . The independent  $x_i$  are  $\pm 1$  and we are dividing their sum by  $N$ . The expected mean of  $A_N$  is still zero. The law of large numbers says that this sample average approaches zero with probability 1. How fast does  $A_N$  approach zero? **What is its variance  $\sigma_N^2$ ?**

$$\text{By linearity } \sigma_N^2 = \frac{\sigma^2}{N^2} + \frac{\sigma^2}{N^2} + \cdots + \frac{\sigma^2}{N^2} = N \frac{\sigma^2}{N^2} = \frac{1}{N} \text{ since } \sigma^2 = 1. \quad (11)$$

**Example 3** Change outputs from 1 or  $-1$  to  $x = 1$  or  $x = 0$ . Keep  $p_1 = p_0 = \frac{1}{2}$ .

The new mean value  $m = \frac{1}{2}$  falls halfway between 0 and 1. The variance moves to  $\sigma^2 = \frac{1}{4}$ :

$$m = \frac{1}{2}(1) + \frac{1}{2}(0) = \frac{1}{2} \quad \text{and} \quad \sigma^2 = \frac{1}{2} \left(1 - \frac{1}{2}\right)^2 + \frac{1}{2} \left(0 - \frac{1}{2}\right)^2 = \frac{1}{4}.$$

The average  $A_N$  now has mean  $\frac{1}{2}$  and variance  $\frac{1}{4N^2} + \cdots + \frac{1}{4N^2} = \frac{1}{4N} = \sigma_N^2$ . (12)

This  $\sigma_N$  is half the size of  $\sigma_N$  in Example 2. This must be correct because the new range 0 to 1 is half as long as  $-1$  to 1. Examples 2-3 are showing a law of linearity.

**The new 0 – 1 variable  $x_{\text{new}}$  is  $\frac{1}{2} x_{\text{old}} + \frac{1}{2}$ .** So the mean  $m$  is increased to  $\frac{1}{2}$  and the variance is *multiplied* by  $\left(\frac{1}{2}\right)^2$ . A shift changes  $m$  and the rescaling changes  $\sigma^2$ .

**Linearity**  $x_{\text{new}} = ax_{\text{old}} + b$  has  $m_{\text{new}} = am_{\text{old}} + b$  and  $\sigma_{\text{new}}^2 = a^2 \sigma_{\text{old}}^2$ . (13)

Here are the results from three numerical tests: random 0 or 1 averaged over  $N$  trials.

[48 1's from  $N = 100$ ] [5035 1's from  $N = 10000$ ] [19967 1's from  $N = 40000$ ].

**The standardized  $X = (x - m)/\sigma = (A_N - \frac{1}{2}) / 2\sqrt{N}$  was  $[-.40]$   $[.70]$   $[-.33]$ .**

The Central Limit Theorem says that the average of many coin flips will approach a normal distribution. Let us begin to see how that happens: **binomial approaches normal**.

For each flip, the probability of heads is  $\frac{1}{2}$ . For  $N = 3$  flips, the probability of heads all three times is  $\left(\frac{1}{2}\right)^3 = \frac{1}{8}$ . The probability of heads twice and tails once is  $\frac{3}{8}$ , from three sequences HHT and HTH and THH. These numbers  $\frac{1}{8}$  and  $\frac{3}{8}$  are pieces of  $\left(\frac{1}{2} + \frac{1}{2}\right)^3 = \frac{1}{8} + \frac{3}{8} + \frac{3}{8} + \frac{1}{8} = 1$ . *The average number of heads in 3 flips is 1.5.*

$$\text{Mean } m = (3 \text{ heads})\frac{1}{8} + (2 \text{ heads})\frac{3}{8} + (1 \text{ head})\frac{3}{8} + 0 = \frac{3}{8} + \frac{6}{8} + \frac{3}{8} = 1.5 \text{ heads}$$



With  $N$  flips, Example 3 (or common sense) gives a mean of  $m = \sum x_i p_i = \frac{1}{2}N$  heads.

The variance  $\sigma^2$  is based on the *squared distance* from this mean  $N/2$ . With  $N = 3$  the variance is  $\sigma^2 = \frac{3}{4}$  (which is  $N/4$ ). To find  $\sigma^2$  we add  $(x_i - m)^2 p_i$  with  $m = 1.5$ :

$$\sigma^2 = (3 - 1.5)^2 \frac{1}{8} + (2 - 1.5)^2 \frac{3}{8} + (1 - 1.5)^2 \frac{3}{8} + (0 - 1.5)^2 \frac{1}{8} = \frac{9 + 3 + 3 + 9}{32} = \frac{3}{4}.$$

For any  $N$ , the variance is  $\sigma_N^2 = N/4$ . Then  $\sigma_N = \sqrt{N}/2$ .

Figure 12.3 shows how the probabilities of 0, 1, 2, 3, 4 heads in  $N = 4$  flips come close to a bell-shaped Gaussian. That Gaussian is centered at the mean value  $N/2 = 2$ . To reach the standard Gaussian (mean 0 and variance 1) we shift and rescale that graph. If  $x$  is the number of heads in  $N$  flips—the average of  $N$  zero-one outcomes—then  $x$  is shifted by its mean  $m = N/2$  and rescaled by  $\sigma = \sqrt{N}/2$  to produce the standard  $X$ :

**Shifted and scaled**

$$X = \frac{x - m}{\sigma} = \frac{x - \frac{1}{2}N}{\sqrt{N}/2} \quad (N = 4 \text{ has } X = x - 2)$$

**Subtracting  $m$  is “centering” or “detrending”. The mean of  $X$  is zero.**

**Dividing by  $\sigma$  is “normalizing” or “standardizing”. The variance of  $X$  is 1.**

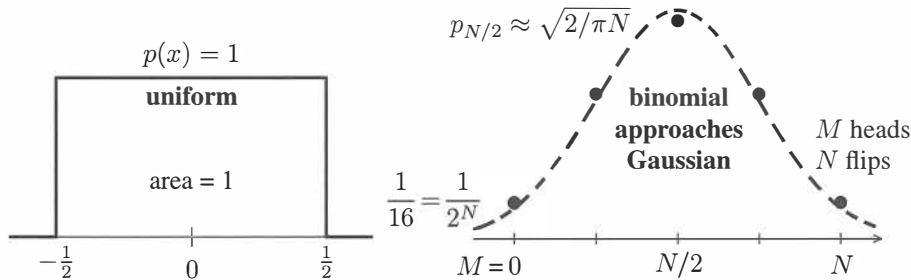


Figure 12.3: The probabilities  $p = (1, 4, 6, 4, 1)/16$  for the number of heads in 4 flips. These  $p_i$  approach a Gaussian distribution with variance  $\sigma^2 = N/4$  centered at  $m = N/2$ . For  $X$ , the Central Limit Theorem gives convergence to the normal distribution  $\mathbf{N}(0, 1)$ .

It is fun to see the Central Limit Theorem giving the right answer at the center point  $X = 0$ . At that point, the factor  $e^{-X^2/2}$  equals 1. We know that the variance for  $N$  coin flips is  $\sigma^2 = N/4$ . The center of the bell-shaped curve has height  $1/\sqrt{2\pi\sigma^2} = \sqrt{2/N\pi}$ .

What is the height at the center of the coin-flip distribution  $p_0$  to  $p_N$  (the binomial distribution)? For  $N = 4$ , the probabilities for 0, 1, 2, 3, 4 heads come from  $(\frac{1}{2} + \frac{1}{2})^4$ .

$$\text{Center probability } \frac{6}{16} \quad \left(\frac{1}{2} + \frac{1}{2}\right)^4 = \frac{1}{16} + \frac{4}{16} + \frac{6}{16} + \frac{4}{16} + \frac{1}{16} = 1.$$



The binomial theorem in Problem 8 tells us the center probability  $p_{N/2}$  for any even  $N$ :

$$\text{The center probability } \left( \frac{N}{2} \text{ heads, } \frac{N}{2} \text{ tails} \right) \text{ is } \frac{1}{2^N} \frac{N!}{(N/2)! (N/2)!}$$

For  $N = 4$ , those factorials produce  $4!/2!2! = 24/4 = 6$ . For large  $N$ , Stirling's formula  $\sqrt{2\pi N}(N/e)^N$  is a close approximation to  $N!$ . Use Stirling for  $N$  and twice for  $N/2$ :

$$\begin{array}{ll} \text{Limit of coin-flip} & \\ \text{Center probability} & p_{N/2} \approx \frac{1}{2^N} \frac{\sqrt{2\pi N}(N/e)^N}{\pi N(N/2e)^N} = \frac{\sqrt{2}}{\sqrt{\pi N}} = \frac{1}{\sqrt{2\pi\sigma}}. \end{array} \quad (14)$$

At that last step we used the variance  $\sigma^2 = N/4$  for the coin-tossing problem. The result  $1/\sqrt{2\pi\sigma}$  matches the center value (above) for the Gaussian. The Central Limit Theorem is true: The “binomial distribution” approaches the normal distribution as  $N \rightarrow \infty$ .

## Monte Carlo Estimation Methods

Scientific computing has to work with errors in the data. Financial computing has to work with unsure numbers and uncertain predictions. All of applied mathematics has moved to **accepting uncertainty in the inputs and estimating the variance in the outputs**.

How to estimate that variance? Often probability distributions  $p(x)$  are not known. What we can do is to try different inputs  $b$  and compute the outputs  $x$  and take an average. This is the simplest form of a **Monte Carlo method** (named after the gambling palace on the Riviera, where I once saw a fight about whether the bet was placed in time). Monte Carlo approximates an expected value  $E[x]$  by a sample average  $(x_1 + \cdots + x_N)/N$ .

Please understand that every  $x_k$  can be expensive to compute. We are not just flipping coins. Each sample comes from a set of data  $b_k$ . *Monte Carlo randomly chooses this data  $b_k$ , it computes the outputs  $x_k$ , and then it averages those  $x$ 's.* Decent accuracy for  $E[x]$  often requires many samples  $b$  and huge computing cost. The error in approximating  $E[x]$  by  $(x_1 + \cdots + x_N)/N$  is normally of order  $1/\sqrt{N}$ . *Slow improvement as  $N$  increases.*

That  $1/\sqrt{N}$  estimate came for coin flips in equation (11). Averaging  $N$  independent samples  $x_k$  of variance  $\sigma^2$  reduces the variance to  $\sigma^2/N$ .

“Quasi-Monte Carlo” can sometimes reduce this variance to  $\sigma^2/N^2$ : a big difference! The inputs  $b_k$  are selected very carefully—not just randomly. This QMC approach is surveyed in the journal *Acta Numerica* 2013. The newer idea of “Multilevel Monte Carlo” is outlined by Michael Giles in *Acta Numerica* 2015. Here is how it works.

Suppose it is much simpler to simulate another variable  $y(b)$  close to  $x(b)$ . Then use  $N$  computations of  $y(b_k)$  and only  $N^* < N$  computations of  $x(b_k)$  to estimate  $E[x]$ .

### 2-level Monte Carlo

$$E[x] \approx \frac{1}{N} \sum_1^N y(b_k) + \frac{1}{N^*} \sum_1^{N^*} [x(b_k) - y(b_k)].$$

The idea is that  $x - y$  has a smaller variance  $\sigma^*$  than the original  $x$ . Therefore  $N^*$  can be smaller than  $N$ , with the same accuracy for  $E[x]$ . We do  $N$  cheap simulations to find the  $y$ 's. Those cost  $C$  each. We only do  $N^*$  expensive simulations involving  $x$ 's. Those cost  $C^*$  each. The total computing cost is  $NC + N^*C^*$ .

Calculus minimizes the overall variance for a fixed total cost. The optimal ratio  $N^*/N$  is  $\sqrt{C/C^*} \sigma^*/\sigma$ . Three-level Monte Carlo would simulate  $x, y$ , and  $z$ :

$$E[x] \approx \frac{1}{N} \sum_1^N z(b_k) + \frac{1}{N^*} \sum_1^{N^*} [y(b_k) - z(b_k)] + \frac{1}{N^{**}} \sum_1^{N^{**}} [x(b_k) - y(b_k)].$$

Giles optimizes  $N, N^*, N^{**}, \dots$  to keep  $E[x] \leq \text{fixed } E_0$ , and provides a MATLAB code.

### Review : Three Formulas for the Mean and the Variance

The formulas for  $m$  and  $\sigma^2$  are the starting point for all of probability and statistics. There are three different cases to keep straight: **sample** values  $X_i$ , **expected** values (discrete  $p_i$ ), and a range of **expected** values (continuous  $p(x)$ ). Here are the mean and the variance:

|  |                                   |   |
|--|-----------------------------------|---|
| <b>Samples <math>X_1</math> to <math>X_N</math></b>                        | $m = \frac{X_1 + \dots + X_N}{N}$ | $S^2 = \frac{(X_1 - m)^2 + \dots + (X_N - m)^2}{N - 1}$ |
| <b><math>n</math> possible outputs with probabilities <math>p_i</math></b> | $m = \sum_1^n p_i x_i$            | $\sigma^2 = \sum_1^n p_i (x_i - m)^2$                   |
| <b>Range of outputs with probability density</b>                           | $m = \int x p(x) dx$              | $\sigma^2 = \int (x - m)^2 p(x) dx$                     |

A natural question: Why are there no probabilities  $p$  on the first line? How can these formulas be parallel? Answer: *We expect a fraction  $p_i$  of the samples to be  $X = x_i$ .* If this is exactly true,  $X = x_i$  is repeated  $p_i N$  times. Then lines 1 and 2 give the same  $m$ .

When we work with samples, we don't know the  $p_i$ . We just include each output  $X$  as often as it comes. We get the "empirical" mean instead of the expected mean.

### Problem Set 12.1

- 1 Add 7 to every output  $x$ . What happens to the mean and the variance? What are the new sample mean, the new expected mean, and the new variance?
- 2 We know:  $\frac{1}{3}$  of all integers are divisible by 3 and  $\frac{1}{7}$  of integers are divisible by 7. What fraction of integers will be divisible by 3 or 7 or both?
- 3 Suppose you sample from the numbers 1 to 1000 with equal probabilities  $1/1000$ . What are the probabilities  $p_0$  to  $p_9$  that the last digit of your sample is  $0, \dots, 9$ ? What is the expected mean  $m$  of that last digit? What is its variance  $\sigma^2$ ?
- 4 Sample again from 1 to 1000 but look at the last digit of the sample *squared*. That square could end with  $x = 0, 1, 4, 5, 6$ , or 9. What are the probabilities  $p_0, p_1, p_4, p_5, p_6, p_9$ ? What are the (expected) mean  $m$  and variance  $\sigma^2$  of that number  $x$ ?

- 5 (a little tricky) Sample again from 1 to 1000 with equal probabilities and let  $x$  be the *first* digit ( $x = 1$  if the number is 15). What are the probabilities  $p_1$  to  $p_9$  (adding to 1) of  $x = 1, \dots, 9$ ? What are the mean and variance of  $x$ ?
- 6 Suppose you have  $N = 4$  samples 157, 312, 696, 602 in Problem 5. What are the first digits  $x_1$  to  $x_4$  of the squares? What is the sample mean  $\mu$ ? What is the sample variance  $S^2$ ? Remember to divide by  $N - 1 = 3$  and not  $N = 4$ .

- 7 Equation (4) gave a second equivalent form for  $S^2$  (the variance using samples):

$$S^2 = \frac{1}{N-1} \text{sum of } (x_i - m)^2 = \frac{1}{N-1} [(\text{sum of } x_i^2) - Nm^2].$$

Verify the matching identity for the expected variance  $\sigma^2$  (using  $m = \sum p_i x_i$ ):

$$\sigma^2 = \text{sum of } p_i (x_i - m)^2 = (\text{sum of } p_i x_i^2) - m^2.$$

- 8 If all 24 samples from a population produce the same age  $x = 20$ , what are the sample mean  $\mu$  and the sample variance  $S^2$ ? What if  $x = 20$  or 21, 12 times each?
- 9 Computer experiment as on page 541: Find the average  $A_{1000000}$  of a million random 0-1 samples! What is  $X = (A_N - \frac{1}{2}) / 2\sqrt{N}$ ?
- 10 The probability  $p_i$  to get  $i$  heads in  $N$  coin flips is the *binomial number*  $b_i = \binom{N}{i}$  divided by  $2^N$ . The  $b_i$  add to  $(1+1)^N = 2^N$  so the probabilities  $p_i$  add to 1.

$$p_0 + \dots + p_N = \left(\frac{1}{2} + \frac{1}{2}\right)^N = \frac{1}{2^N} (b_0 + \dots + b_N) \text{ with } b_i = \frac{N!}{i!(N-i)!}$$

$$N=4 \text{ leads to } b_0 = \frac{24}{24}, b_1 = \frac{24}{(1)(6)} = 4, b_2 = \frac{24}{(2)(2)} = 6, p_i = \frac{1}{16} (1, 4, 6, 4, 1).$$

Notice  $b_i = b_{N-i}$ . *Problem:* Confirm that the mean  $m = 0p_0 + \dots + Np_N$  equals  $\frac{N}{2}$ .

- 11 For any function  $f(x)$  the expected value is  $E[f] = \sum p_i f(x_i)$  or  $\int p(x) f(x) dx$  (discrete probability or continuous probability). Suppose the mean is  $E[x] = m$  and the variance is  $E[(x - m)^2] = \sigma^2$ . **What is  $E[x^2]$ ?**
- 12 Show that the standard normal distribution  $p(x)$  has total probability  $\int p(x) dx = 1$  as required. A famous trick multiplies  $\int p(x) dx$  by  $\int p(y) dy$  and computes the integral over all  $x$  and all  $y$  ( $-\infty$  to  $\infty$ ). The trick is to replace  $dx dy$  in that double integral by  $r dr d\theta$  (polar coordinates with  $x^2 + y^2 = r^2$ ). Explain each step:

$$2\pi \int_{-\infty}^{\infty} p(x) dx \int_{-\infty}^{\infty} p(y) dy = \iiint_{-\infty}^{\infty} e^{-(x^2+y^2)/2} dx dy = \int_{\theta=0}^{2\pi} \int_{r=0}^{\infty} e^{-r^2/2} r dr d\theta = 2\pi.$$

## 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  $\mathbf{m} = (m_1, m_2, m_3)$  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  $\sigma_i^2 = E[(x_i - m_i)^2]$  based on the squared distance from its mean. Those  $M$  numbers  $\sigma_1^2, \dots, \sigma_M^2$  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 ( $a, h, w$ ) 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.**

|   |     |
|---|-----|
| <b>Covariance</b> $\sigma_{ah} = E[(\text{age} - \text{mean age})(\text{height} - \text{mean height})]$ . | (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.

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

$p_{ij}$  = **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  $y_j$ :

|  |     |
|--|-----|
| <b>Covariance</b> $\sigma_{12} = \sum_{\text{all } i, j} p_{ij}(x_i - m_1)(y_j - m_2)$ | (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  $(i, j) = (\text{Prob of } i)(\text{Prob of } j)$ .**

**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**

$$P = \begin{bmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{bmatrix} = \begin{bmatrix} \frac{1}{4} & \frac{1}{4} \\ \frac{1}{4} & \frac{1}{4} \end{bmatrix} \quad P = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{bmatrix}.$$

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  $(x_1, y_1) = (\text{heads}, \text{heads})$  and  $(x_1, y_2) = (\text{heads}, \text{tails})$ . Notice the row sums  $p_i$  and column sums  $P_j$  and the total sum = 1.

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

(second coin) column sums  $P_1 \quad P_2$       4 entries add to 1

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

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

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

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

**Zero covariance  $\sigma_{12}$  for independent trials**       $V = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} = \text{diagonal covariance matrix.}$

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

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

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  $\sigma_{12}$ :

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

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

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

**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  $X = (\text{age}, \text{height}, \text{weight})$  contributes to the **sample mean  $\bar{X}$**  and to  $S$ . Each term  $(X_i - \bar{X})(X_i - \bar{X})^T$  is positive semidefinite and we just add to reach  $S$ :

$$\bar{\mathbf{X}} = \frac{\mathbf{X}_1 \oplus \cdots \oplus \mathbf{X}_N}{N} \quad \mathbf{S} = \frac{(\mathbf{X}_1 - \bar{\mathbf{X}})(\mathbf{X}_1 - \bar{\mathbf{X}})^T + \cdots + (\mathbf{X}_N - \bar{\mathbf{X}})(\mathbf{X}_N - \bar{\mathbf{X}})^T}{N - 1} \quad (3)$$

### The Covariance Matrix $V$ is Positive Semidefinite

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

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

$p_{ij} \geq 0$  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.

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

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  $(x_i, y_1), (x_i, y_2), \dots, (x_i, y_n)$ :

$$\text{Row sum } p_i \text{ of } P \quad \sum_{j=1}^n p_{ij} = \text{probability } p_i \text{ of } x_i \text{ in experiment 1.} \quad (6)$$

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:

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

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  $(i, j)$ :

$$\begin{aligned} \text{Covariance matrix} \quad V &= \sum \sum V_{ij} \quad V = \sum_{\text{all } i, j} p_{ij} \begin{bmatrix} (x_i - m_1)^2 & (x_i - m_1)(y_j - m_2) \\ (x_i - m_1)(y_j - m_2) & (y_j - m_2)^2 \end{bmatrix} \end{aligned} \quad (7)$$

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:

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

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  $U^T$ :

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

Every matrix  $UU^T$  is positive semidefinite. So the whole matrix  $V$  (combining these matrices  $UU^T$  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  $\bar{X} = E[X]$ :

$$\text{Covariance matrix} \quad V = E \left[ (X - \bar{X})(X - \bar{X})^T \right] \quad (10)$$

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

For  $M = 1$  (one variable) you see that  $\bar{X}$  is the mean  $m$  and  $V$  is  $\sigma^2$  (Section 12.1). For  $M = 2$  (two coins) you see that  $\bar{X}$  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  $X = (5 \text{ years}, 31 \text{ inches}, 48 \text{ pounds})$  and its probability is  $p_{5,31,48}$ .

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

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

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

$$\begin{aligned} \text{Variance of } c^T X &= E \left[ (c^T X - c^T \bar{X})(c^T X - c^T \bar{X})^T \right] \\ &= c^T E \left[ (X - \bar{X})(X - \bar{X})^T \right] c = c^T V c! \end{aligned} \quad (11)$$

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

Covariance matrices  $V$  open up the link between probability and linear algebra:  $V$  equals  $Q\Lambda Q^T$  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  $\mathbf{c}^T \mathbf{V} \mathbf{c} \geq 0$ . 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  $UU^T$ . Here  $U$  is  $X - \bar{X}$ :

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

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

The value of the expectation symbol  $\mathbf{E}$  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  $i = 0, 1, 2, 3, \dots$ , 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  $V = \mathbf{E}[UU^T]$ :

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

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

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

**Dependent variables**  $x, y, z$      $p(x, y, z) = 0$  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  $m_z = m_x + m_y$ :**

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

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*:

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

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  $\mathbf{E}[x]$ :

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

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  $\dots$  and  $(x_N, y_j)$  are all the ways to produce  $y_j$ :

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

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

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:

$$\begin{aligned} \sigma_z^2 &= \sum_i \sum_j p_{ij} (x_i + y_j - m_x - m_y)^2 \\ &= \sum_i \sum_j p_{ij} (x_i - m_x)^2 + \sum_i \sum_j p_{ij} (y_j - m_y)^2 + 2 \sum_i \sum_j p_{ij} (x_i - m_x)(y_j - m_y) \end{aligned}$$

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

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

### The Covariance Matrix for $\mathbf{Z} = \mathbf{AX}$

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

$$\sigma_z^2 = \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_y^2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \text{which is } \sigma_z^2 = \mathbf{AVAT}. \quad (17)$$

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

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

That vector  $\mathbf{Z} = \mathbf{AX}$  of length  $K$  has a  $K$  by  $K$  covariance matrix  $\mathbf{V}_Z$ . 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  $\mathbf{AX}$  is  $A$  (covariance matrix of  $\mathbf{X}$ )  $A^T$ :

$$\boxed{\text{The covariance matrix of } \mathbf{Z} = \mathbf{AX} \text{ is } \mathbf{V}_Z = \mathbf{AV}_X \mathbf{A}^T} \quad (18)$$

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  $\sigma_y$ :

|   |  |
|---|--|
| <b>Correlation</b> $\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} = \text{covariance of } \frac{x}{\sigma_x} \text{ and } \frac{y}{\sigma_y}$ | <b>Always</b> $-1 \leq \rho_{xy} \leq 1$ |
|---|--|

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

We know that always  $\sigma_{xy}^2 \leq \sigma_x^2 \sigma_y^2$  (the covariance matrix  $V$  is at least positive semidefinite). Then  $\rho_{xy}^2 \leq 1$ . 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):

|  |  |
|--|--|
| <b>Correlation matrix</b> $R = \begin{bmatrix} 1 & \rho_{xy} \\ \rho_{xy} & 1 \end{bmatrix}$ | $R = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$ 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  $D = \text{diag}[1/\sigma_1, \dots, 1/\sigma_M]$ .
- (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**

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

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:

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

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  $ax + by$ :

$$\text{Use } V_Z = AV_X A^T \quad \sigma_{ax+by}^2 = \sigma_{ax}^2 + \sigma_{by}^2 + 2\sigma_{ax,by} = a^2 + b^2 + 0.$$

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:

$$\begin{array}{ll} \text{Difference matrix} & \text{Covariance matrix} \\ A = \begin{bmatrix} 1 & -1 \end{bmatrix} & V_Z = AV_X A^T \end{array} \quad V_Z = \begin{bmatrix} 1 & -1 \end{bmatrix} \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} \\ & = \sigma_1^2 - 2\sigma_{12} + \sigma_2^2$$

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  $\sigma^2 = \underline{\hspace{2cm}}$ .
- 2 What is the covariance  $\sigma_{kl}$  between the results  $x_1, \dots, x_n$  of Experiment 3 and the results  $y_1, \dots, y_n$  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**( $Y = \text{heads}$  when  $X = \text{heads}$ ) is 1.

*Example 2* Independent coin flips:  $X$  gives no information about  $Y$ . Useless to know  $X$ .  
 Then **Prob**( $Y = \text{heads} | X = \text{heads}$ ) is the same as **Prob**( $Y = \text{heads}$ ).

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

$$\text{Prob}(Y = y_j) = \text{sum over all outputs } x_i \text{ of } \text{Prob}(Y = y_j | X = x_i).$$
- 6 The  $n$  by  $n$  matrix  $P$  contains **joint probabilities**  $p_{ij} = \text{Prob}(X = x_i \text{ and } Y = y_j)$ .  
 Explain why the conditional  $\text{Prob}(Y = y_j | X = x_i)$  equals  $\frac{p_{ij}}{p_{i1} + \dots + p_{in}} = \frac{p_{ij}}{p_i}$ .
- 7 For this joint probability matrix with  $\text{Prob}(x_1, y_2) = 0.3$ , find  $\text{Prob}(y_2 | x_1)$  and  $\text{Prob}(x_1)$ .  

$$P = \begin{bmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{bmatrix} = \begin{bmatrix} 0.1 & 0.3 \\ 0.2 & 0.4 \end{bmatrix} \quad \begin{array}{l} \text{The entries } p_{ij} \text{ add to 1.} \\ \text{Some } i, j \text{ must happen.} \end{array}$$
- 8 Explain the **product rule** of conditional probability:  
 $p_{ij} = \text{Prob}(X = x_i \text{ and } Y = y_j)$  equals  $\text{Prob}(Y = y_j | X = x_i)$  times  $\text{Prob}(X = x_i)$ .
- 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.

## 12.3 Multivariate Gaussian and Weighted Least Squares

The normal probability density  $p(x)$  (the Gaussian) depends on only two numbers :

**Mean  $m$  and variance  $\sigma^2$**  
$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-m)^2/2\sigma^2}. \quad (1)$$

The graph of  $p(x)$  is a bell-shaped curve centered at  $x = m$ . The continuous variable  $x$  can be anywhere between  $-\infty$  and  $\infty$ . With probability close to  $\frac{2}{3}$ , that random  $x$  will lie between  $m - \sigma$  and  $m + \sigma$  (less than one standard deviation  $\sigma$  from its mean value  $m$ ).

$$\int_{-\infty}^{\infty} p(x) dx = 1 \quad \text{and} \quad \int_{m-\sigma}^{m+\sigma} p(x) dx = \frac{1}{\sqrt{2\pi}} \int_{-1}^1 e^{-X^2/2} dX \approx \frac{2}{3}. \quad (2)$$

That integral has a change of variables from  $x$  to  $X = (x - m)/\sigma$ . This simplifies the exponent to  $-X^2/2$  and it simplifies the limits of integration to  $-1$  and  $1$ . Even the  $1/\sigma$  from  $p$  disappears outside the integral because  $dX$  equals  $dx/\sigma$ . Every Gaussian turns into a **standard Gaussian**  $p(X)$  with mean  $m = 0$  and variance  $\sigma^2 = 1$ . Just call it  $p(x)$  :

**The standard normal distribution  $N(0, 1)$**  has 
$$p(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}. \quad (3)$$

Integrating  $p(x)$  from  $-\infty$  to  $x$  gives the cumulative distribution  $F(x)$ : the probability that a random sample is below  $x$ . That probability will be  $F = \frac{1}{2}$  at  $x = 0$  (the mean).

### Two-dimensional Gaussians

Now we have  $M = 2$  Gaussian random variables  $x$  and  $y$ . They have means  $m_1$  and  $m_2$ . They have variances  $\sigma_1^2$  and  $\sigma_2^2$ . If they are *independent*, then their probability density  $p(x, y)$  is just  $p_1(x)$  **times**  $p_2(y)$ . Multiply probabilities when variables are independent :

**Independent  $x$  and  $y$**  
$$p(x, y) = \frac{1}{2\pi\sigma_1\sigma_2} e^{-(x-m_1)^2/2\sigma_1^2} e^{-(y-m_2)^2/2\sigma_2^2} \quad (4)$$

The covariance of  $x$  and  $y$  will be  $\sigma_{12} = 0$ . The covariance matrix  $V$  will be *diagonal*. The variances  $\sigma_1^2$  and  $\sigma_2^2$  are always on the main diagonal of  $V$ . The exponent in  $p(x, y)$  is just the sum of the  $x$ -exponent and the  $y$ -exponent. Good to notice that the two exponents can be combined into  $-\frac{1}{2}(x-m)^T V^{-1}(x-m)$  with  $V^{-1}$  in the middle :

$$-\frac{(x-m_1)^2}{2\sigma_1^2} - \frac{(y-m_2)^2}{2\sigma_2^2} = -\frac{1}{2} \begin{bmatrix} x-m_1 & y-m_2 \end{bmatrix} \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}^{-1} \begin{bmatrix} x-m_1 \\ y-m_2 \end{bmatrix} \quad (5)$$



### Non-independent $x$ and $y$

We are ready to give up independence. The exponent (5) with  $V^{-1}$  is still correct when  $V$  is no longer a diagonal matrix. **Now the Gaussian depends on a vector  $m$  and a matrix  $V$ .**

When  $M = 2$ , the first variable  $x$  may give partial information about the second variable  $y$  (and vice versa). Maybe part of  $y$  is decided by  $x$  and part is truly independent. It is the  $M$  by  $M$  covariance matrix  $V$  that accounts for dependencies between the  $M$  variables  $x = x_1, \dots, x_M$ . Its inverse  $V^{-1}$  goes into  $p(x)$ :

**Multivariate Gaussian probability distribution**

$$p(x) = \frac{1}{(\sqrt{2\pi})^M \sqrt{\det V}} e^{-(x-m)^T V^{-1} (x-m)/2} \quad (6)$$

The vectors  $x = (x_1, \dots, x_M)$  and  $m = (m_1, \dots, m_M)$  contain the random variables and their means. The  $M$  square roots of  $2\pi$  and the determinant of  $V$  are included to make the total probability equal to 1. Let me check that by linear algebra. I use the eigenvalues  $\lambda$  and orthonormal eigenvectors  $q$  of the symmetric matrix  $V = Q\Lambda Q^T$ . So  $V^{-1} = Q\Lambda^{-1}Q^T$ :

$$X = x - m \quad (x - m)^T V^{-1} (x - m) = X^T Q \Lambda^{-1} Q^T X = Y^T \Lambda^{-1} Y$$

*Notice!* The combinations  $Y = Q^T X = Q^T (x - m)$  are statistically independent. *Their covariance matrix  $\Lambda$  is diagonal.*

This step of diagonalizing  $V$  by its eigenvector matrix  $Q$  is the same as “uncorrelating” the random variables. Covariances are zero for the new variables  $X_1, \dots, X_M$ . This is the point where linear algebra helps calculus to compute multidimensional integrals.

The integral of  $p(x)$  is not changed when we center the variable  $x$  by subtracting  $m$  to reach  $X$ , and rotate that variable to reach  $Y = Q^T X$ . The matrix  $\Lambda$  is diagonal! So the integral we want splits into  $M$  separate one-dimensional integrals that we know:

$$\begin{aligned} \int \dots \int e^{-Y^T \Lambda^{-1} Y/2} dY &= \left( \int_{-\infty}^{\infty} e^{-y_1^2/2\lambda_1} dy_1 \right) \dots \left( \int_{-\infty}^{\infty} e^{-y_M^2/2\lambda_M} dy_M \right) \\ &= \left( \sqrt{2\pi\lambda_1} \right) \dots \left( \sqrt{2\pi\lambda_M} \right) = \left( \sqrt{2\pi} \right)^M \sqrt{\det V}. \end{aligned} \quad (7)$$

The determinant of  $V$  (also the determinant of  $\Lambda$ ) is the product  $(\lambda_1) \dots (\lambda_M)$  of the eigenvalues. Then (7) gives the correct number to divide by so that  $p(x_1, \dots, x_M)$  in equation (6) has integral = 1 as desired.

The mean and variance of  $p(x)$  are also  $M$ -dimensional integrals. The same idea of diagonalizing  $V$  by its eigenvectors and introducing  $Y = Q^T X$  will find those integrals:

$$\text{Vector } m \text{ of means} \quad \int \dots \int x p(x) dx = (m_1, m_2, \dots) = m \quad (8)$$

$$\text{Covariance matrix } V \quad \int \dots \int (x - m) p(x) (x - m)^T dx = V. \quad (9)$$

Conclusion: Formula (6) for the probability density  $p(x)$  has all the properties we want.



## Weighted Least Squares

In Chapter 4, least squares started from an unsolvable system  $A\mathbf{x} = \mathbf{b}$ . We chose  $\hat{\mathbf{x}}$  to minimize the error  $\|\mathbf{b} - A\mathbf{x}\|^2$ . That led us to the least squares equation  $A^T A \hat{\mathbf{x}} = A^T \mathbf{b}$ . The best  $A\hat{\mathbf{x}}$  is the projection of  $\mathbf{b}$  onto the column space of  $A$ . But is this squared distance  $E = \|\mathbf{b} - A\mathbf{x}\|^2$  the right error measure to minimize?

If the measurement errors in  $\mathbf{b}$  are independent random variables, with mean  $m = 0$  and variance  $\sigma^2 = 1$  and a normal distribution, Gauss would say **yes**: *Use least squares*. If the errors are not independent or their variances are not equal. Gauss would say **no**: *Use **weighted** least squares*. This section will show that the good measure of error is  $E = (\mathbf{b} - A\mathbf{x})^T V^{-1} (\mathbf{b} - A\mathbf{x})$ . The equation for the best  $\hat{\mathbf{x}}$  uses the covariance matrix  $V$ :

**Weighted least squares**

$$A^T V^{-1} A \hat{\mathbf{x}} = A^T V^{-1} \mathbf{b}. \quad (10)$$

The most important examples have  $m$  *independent* errors in  $\mathbf{b}$ . Those errors have variances  $\sigma_1^2, \dots, \sigma_m^2$ . By independence,  $V$  is a diagonal matrix. The good weights  $1/\sigma_1^2, \dots, 1/\sigma_m^2$  come from  $V^{-1}$ . *We are weighting the errors in  $\mathbf{b}$  to have **variance** = 1:*

**Weighted least squares  
Independent errors in  $\mathbf{b}$**

$$\text{Minimize } E = \sum_{i=1}^m \frac{(\mathbf{b} - A\mathbf{x})_i^2}{\sigma_i^2} \quad (11)$$

By weighting the errors, we are “whitening” the noise. **White noise** is a quick description of independent errors based on the standard Gaussian  $\mathbf{N}(0, 1)$  with mean zero and  $\sigma^2 = 1$ .

Let me write down the steps to equations (10) and (11) for the best  $\hat{\mathbf{x}}$ :

Start with  $A\mathbf{x} = \mathbf{b}$  ( $m$  equations,  $n$  unknowns,  $m > n$ , no solution)

Each right side  $b_i$  has mean zero and variance  $\sigma_i^2$ . The  $b_i$  are independent.

Divide the  $i$ th equation by  $\sigma_i$  to have variance = 1 for every  $b_i/\sigma_i$

That division turns  $A\mathbf{x} = \mathbf{b}$  into  $V^{-1/2} A \mathbf{x} = V^{-1/2} \mathbf{b}$  with  $V^{-1/2} = \text{diag}(1/\sigma_1, \dots, 1/\sigma_m)$

Ordinary least squares on those weighted equations has  $A \rightarrow V^{-1/2} A$  and  $\mathbf{b} \rightarrow V^{-1/2} \mathbf{b}$

$$(V^{-1/2} A)^T (V^{-1/2} A) \hat{\mathbf{x}} = (V^{-1/2} A)^T V^{-1/2} \mathbf{b} \text{ is } A^T V^{-1} A \hat{\mathbf{x}} = A^T V^{-1} \mathbf{b}. \quad (12)$$

Because of  $1/\sigma^2$  in  $V^{-1}$ , more reliable equations (*smaller*  $\sigma$ ) get heavier weights. This is the main point of weighted least squares.

Those diagonal weightings (uncoupled equations) are the most frequent and the simplest. They apply to *independent errors in the  $b_i$* . When these measurement errors are not independent,  $V$  is no longer diagonal—but (12) is still the correct weighted equation.

In practice, finding all the covariances can be serious work. Diagonal  $V$  is simpler.

### The Variance in the Estimated $\hat{x}$

One more point: Often the important question is not the best  $\hat{x}$  for one particular set of measurements  $\mathbf{b}$ . This is only one sample! The real goal is to know the reliability of the whole experiment. That is measured (as reliability always is) by the **variance in the estimate  $\hat{x}$** . First, zero mean in  $\mathbf{b}$  gives zero mean in  $\hat{x}$ . Then the formula connecting variance  $\mathbf{V}$  in the inputs  $\mathbf{b}$  to variance  $\mathbf{W}$  in the outputs  $\hat{x}$  turns out to be beautiful:

$$\text{Variance-covariance matrix } \mathbf{W} \text{ for } \hat{x} \quad E[(\hat{x} - x)(\hat{x} - x)^T] = (\mathbf{A}^T \mathbf{V}^{-1} \mathbf{A})^{-1}. \quad (13)$$

That smallest possible variance comes from the best possible weighting, which is  $\mathbf{V}^{-1}$ .

This key formula is a perfect application of Section 12.2. **If  $\mathbf{b}$  has covariance matrix  $\mathbf{V}$ , then  $\hat{x} = \mathbf{L}\mathbf{b}$  has covariance matrix  $\mathbf{L}\mathbf{V}\mathbf{L}^T$ .** Equation (12) above tells us that  $\mathbf{L}$  is  $(\mathbf{A}^T \mathbf{V}^{-1} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{V}^{-1}$ . Now substitute this into  $\mathbf{L}\mathbf{V}\mathbf{L}^T$  and watch equation (13) appear:

$$\mathbf{L}\mathbf{V}\mathbf{L}^T = (\mathbf{A}^T \mathbf{V}^{-1} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{V}^{-1} \mathbf{V} \mathbf{V}^{-1} \mathbf{A} (\mathbf{A}^T \mathbf{V}^{-1} \mathbf{A})^{-1} = (\mathbf{A}^T \mathbf{V}^{-1} \mathbf{A})^{-1}.$$

This is the covariance  $\mathbf{W}$  of the output, our best estimate  $\hat{x}$ . It is time for examples.

**Example 1** Suppose a doctor measures your heart rate  $x$  three times ( $m = 3, n = 1$ ):

$$\begin{array}{l} x = b_1 \\ x = b_2 \\ x = b_3 \end{array} \quad \text{is} \quad \mathbf{A}\mathbf{x} = \mathbf{b} \quad \text{with} \quad \mathbf{A} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \quad \text{and} \quad \mathbf{V} = \begin{bmatrix} \sigma_1^2 & 0 & 0 \\ 0 & \sigma_2^2 & 0 \\ 0 & 0 & \sigma_3^2 \end{bmatrix}$$

The variances could be  $\sigma_1^2 = 1/9$  and  $\sigma_2^2 = 1/4$  and  $\sigma_3^2 = 1$ . You are getting more nervous as measurements are taken:  $b_3$  is less reliable than  $b_2$  and  $b_1$ . All three measurements contain some information, so they all go into the best (weighted) estimate  $\hat{x}$ :

$$\begin{array}{l} 3x = 3b_1 \\ 2x = 2b_2 \\ 1x = 1b_3 \end{array} \quad \text{leading to} \quad \mathbf{A}^T \mathbf{V}^{-1} \mathbf{A} \hat{x} = \mathbf{A}^T \mathbf{V}^{-1} \mathbf{b}$$

$$\begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 9 & & \\ & 4 & \\ & & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \hat{x} = \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 9 & & \\ & 4 & \\ & & 1 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

$$\hat{x} = \frac{9b_1 + 4b_2 + b_3}{14} \quad \text{is a weighted average of } b_1, b_2, b_3$$

Most weight is on  $b_1$  since its variance  $\sigma_1$  is smallest. The variance of  $\hat{x}$  has the beautiful formula  $W = (A^T V^{-1} A)^{-1} = 1/14$ :

$$\text{Variance of } \hat{x} \quad \left( \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 9 & & \\ & 4 & \\ & & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right)^{-1} = \frac{1}{14} \quad \text{is smaller than} \quad \frac{1}{9}$$

The **BLUE** theorem of Gauss (proved on the website) says that our  $\hat{x} = Lb$  is the best linear unbiased estimate of the solution to  $Ax = b$ . Any other unbiased choice  $x^* = L^*b$  has greater variance than  $\hat{x}$ . All unbiased choices have  $L^*A = I$  so that an exact  $Ax = b$  will produce the right answer  $x = L^*b = L^*Ax$ .

*Note.* I must add that there are reasons not to minimize squared errors in the first place. One reason: This  $\hat{x}$  often has many small components. The squares of small numbers are very small, and they appear when we minimize. It is easier to make sense of *sparse* vectors—only a few nonzeros. Statisticians often prefer to minimize **unsquared errors**: **the sum of**  $|(b - Ax)_i|$ . *This error measure is  $L^1$  instead of  $L^2$ .* Because of the absolute values, the equation for  $\hat{x}$  becomes nonlinear (it is actually piecewise linear).

Fast new algorithms are computing a sparse  $\hat{x}$  quickly and the future may belong to  $L^1$ .

### The Kalman Filter

The “Kalman filter” is the great algorithm in dynamic least squares. That word *dynamic* means that new measurements  $b_k$  keep coming. So the best estimate  $\hat{x}_k$  keeps changing (based on all of  $b_0, \dots, b_k$ ). More than that, the matrix  $A$  is also changing. So  $\hat{x}_2$  will be our best least squares estimate of the latest solution  $x_k$  to the **whole history of observation equations and update equations (state equations) up to time 2**:

$$A_0 x_0 = b_0 \quad x_1 = F_0 x_0 \quad A_1 x_1 = b_1 \quad x_2 = F_1 x_1 \quad A_2 x_2 = b_2 \quad (14)$$

The Kalman idea is to introduce one equation at a time. There will be errors in each equation. With every new equation, we update the best estimate  $\hat{x}_k$  for the current  $x_k$ . But history is not forgotten! This new estimate  $\hat{x}_k$  uses all the past observations  $b_0$  to  $b_{k-1}$  and all the state equations  $x_{\text{new}} = F_{\text{old}} x_{\text{old}}$ . A large and growing least squares problem.

One more important point. Each least squares equation is **weighted** using the covariance matrix  $V_k$  for the error in  $b_k$ . There is even a covariance matrix  $C_k$  for errors in the update equations  $x_{k+1} = F_k x_k$ . The best  $\hat{x}_2$  then depends on  $b_0, b_1, b_2$  and  $V_0, V_1, V_2$  and  $C_1, C_2$ . The good way to write  $\hat{x}_k$  is as an update to the previous  $\hat{x}_{k-1}$ .

Let me concentrate on a simplified problem, without the matrices  $F_k$  and the covariances  $C_k$ . We are estimating the same true  $x$  at every step. How do we get  $\hat{x}_1$  from  $\hat{x}_0$ ?

$$\text{OLD} \quad A_0 x_0 = b_0 \text{ leads to the weighted equation } A_0^T V_0^{-1} A_0 \hat{x}_0 = A_0^T V_0^{-1} b_0. \quad (15)$$

$$\text{NEW} \quad \begin{bmatrix} A_0 \\ A_1 \end{bmatrix} \hat{x}_1 = \begin{bmatrix} b_0 \\ b_1 \end{bmatrix} \text{ leads to the following weighted equation for } \hat{x}_1 :$$

$$\begin{bmatrix} A_0^T & A_1^T \end{bmatrix} \begin{bmatrix} V_0^{-1} \\ V_1^{-1} \end{bmatrix} \begin{bmatrix} A_0 \\ A_1 \end{bmatrix} \hat{x}_1 = \begin{bmatrix} A_0^T & A_1^T \end{bmatrix} \begin{bmatrix} V_0^{-1} \\ V_1^{-1} \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \end{bmatrix}. \quad (16)$$

Yes, we could just solve that new problem and forget the old one. But the old solution  $\hat{x}_0$  needed work that we hope to reuse in  $\hat{x}_1$ . What we look for is **an update to  $\hat{x}_0$** :

|   |  |
|---|--|
| <b>Kalman update gives <math>\hat{x}_1</math> from <math>\hat{x}_0</math></b> | $\hat{x}_1 = \hat{x}_0 + K_1(b_1 - A_1 \hat{x}_0).$ (17) |
|---|--|

The update correction is the mismatch  $b_1 - A_1 \hat{x}_0$  between the old state  $\hat{x}_0$  and the new measurements  $b_1$ —multiplied by the *Kalman gain matrix*  $K_1$ . The formula for  $K_1$  comes from comparing the solutions  $\hat{x}_1$  and  $\hat{x}_0$  to (15) and (16). And when we update  $\hat{x}_0$  to  $\hat{x}_1$  based on new data  $b_1$ , **we also update the covariance matrix  $W_0$  to  $W_1$** . Remember  $W_0 = (A_0^T V_0^{-1} A_0)^{-1}$  from equation (13). Update its inverse to  $W_1^{-1}$ :

|  |   |
|--|---|
| <b>Covariance <math>W_1</math> of errors in <math>\hat{x}_1</math></b> | $W_1^{-1} = W_0^{-1} + A_1^T V_1^{-1} A_1$ (18) |
|--|---|

|  |                                 |
|--|---------------------------------|
| <b>Kalman gain matrix <math>K_1</math></b> | $K_1 = W_1 A_1^T V_1^{-1}$ (19) |
|--|---------------------------------|

This is the heart of the Kalman filter. Notice the importance of the  $W_k$ . Those matrices measure the reliability of the whole process, where the vector  $\hat{x}_k$  estimates the current state based on the particular measurements  $b_0$  to  $b_k$ .

Whole chapters and whole books are written to explain the dynamic Kalman filter, when the states  $x_k$  are also changing (based on the matrices  $F_k$ ). There is a *prediction* of  $x_k$  using  $F$ , followed by a *correction* using the new data  $b$ . Perhaps best to stop here.

This page was about **recursive least squares**: adding new data  $b_k$  and updating both  $\hat{x}$  and  $W$ : the best current estimate based on all the data, and its covariance matrix.

### Problem Set 12.3

- 1 Two measurements of the same variable  $x$  give two equations  $x = b_1$  and  $x = b_2$ . Suppose the means are zero and the variances are  $\sigma_1^2$  and  $\sigma_2^2$ , with independent errors:  $V$  is diagonal with entries  $\sigma_1^2$  and  $\sigma_2^2$ . Write the two equations as  $Ax = b$  ( $A$  is 2 by 1). As in the text Example 1, find this best estimate  $\hat{x}$  based on  $b_1$  and  $b_2$ :

$$\hat{x} = \frac{b_1/\sigma_1^2 + b_2/\sigma_2^2}{1/\sigma_1^2 + 1/\sigma_2^2} \quad E[\hat{x} \hat{x}^T] = \left( \frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2} \right)^{-1}.$$

- 2 (a) In Problem 1, suppose the second measurement  $b_2$  becomes super-exact and its variance  $\sigma_2 \rightarrow 0$ . What is the best estimate  $\hat{x}$  when  $\sigma_2$  reaches zero?
- (b) The opposite case has  $\sigma_2 \rightarrow \infty$  and no information in  $b_2$ . What is now the best estimate  $\hat{x}$  based on  $b_1$  and  $b_2$ ?

- 3 If  $x$  and  $y$  are independent with probabilities  $p_1(x)$  and  $p_2(y)$ , then  $p(x, y) = p_1(x)p_2(y)$ . By separating double integrals into products of single integrals ( $-\infty$  to  $\infty$ ) show that

$$\iint p(x, y) dx dy = \mathbf{1} \quad \text{and} \quad \iint (x + y) p(x, y) dx dy = \mathbf{m}_1 + \mathbf{m}_2.$$

- 4 Continue Problem 3 for independent  $x, y$  to show that  $p(x, y) = p_1(x)p_2(y)$  has

$$\iint (x - m_1)^2 p(x, y) dx dy = \sigma_1^2 \quad \iint (x - m_1)(y - m_2) p(x, y) dx dy = \mathbf{0}.$$

So the 2 by 2 covariance matrix  $V$  is diagonal and its entries are \_\_\_\_.

- 5 Show that the inverse of a 2 by 2 covariance matrix  $V$  is

$$V^{-1} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix}^{-1} = \frac{\mathbf{1}}{1 - \rho^2} \begin{bmatrix} 1/\sigma_1^2 & -\rho/\sigma_1\sigma_2 \\ -\rho/\sigma_1\sigma_2 & 1/\sigma_2^2 \end{bmatrix} \quad \begin{array}{l} \text{with correlation} \\ \rho = \sigma_{12}/\sigma_1\sigma_2. \end{array}$$

This produces the exponent  $-(x - m)^T V^{-1}(x - m)$  in a 2-variable Gaussian.

- 6 Suppose  $\hat{x}_k$  is the average of  $b_1, \dots, b_k$ . A new measurement  $b_{k+1}$  arrives and we want the new average  $\hat{x}_{k+1}$ . The Kalman update equation (17) is

$$\text{New average} \quad \hat{x}_{k+1} = \hat{x}_k + \frac{1}{k+1} (b_{k+1} - \hat{x}_k).$$

Verify that  $\hat{x}_{k+1}$  is the correct average of  $b_1, \dots, b_{k+1}$ .

- 7 Also check the update equation (18) for the variance  $W_{k+1} = \sigma^2/(k+1)$  of this average  $\hat{x}$  assuming that  $W_k = \sigma^2/k$  and  $b_{k+1}$  has variance  $V = \sigma^2$ .
- 8 (**Steady model**) Problems 6–7 were *static* least squares. All the sample averages  $\hat{x}_k$  were estimates of the same  $x$ . To make the Kalman filter *dynamic*, include also a *state equation*  $x_{k+1} = Fx_k$  with its own error variance  $s^2$ . The dynamic least squares problem allows  $x$  to “drift” as  $k$  increases:

$$\begin{bmatrix} 1 & \\ -F & 1 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} b_0 \\ b_1 \end{bmatrix} \quad \text{with variances} \quad \begin{bmatrix} \sigma^2 \\ s^2 \\ \sigma^2 \end{bmatrix}.$$

With  $F = 1$ , divide both sides of those three equations by  $\sigma, s$ , and  $\sigma$ . Find  $\hat{x}_0$  and  $\hat{x}_1$  by least squares, which gives more weight to the recent  $b_1$ . The Kalman filter is developed in *Algorithms for Global Positioning* (Borre and Strang, Wellesley-Cambridge Press).

### Change in $A^{-1}$ from a Change in $A$

This final page connects the beginning of the book (inverses and rank one matrices) with the end of the book (dynamic least squares and filters). Begin with this basic formula:

**The inverse of  $M = I - uv^T$  is  $M^{-1} = I + \frac{uv^T}{1 - v^T u}$**

The quickest proof is  $MM^{-1} = I - uv^T + (1 - uv^T) \frac{uv^T}{1 - v^T u} = I - uv^T + uv^T = I$ .

$M$  is not invertible if  $v^T u = 1$  (then  $Mu = 0$ ). Here  $v^T = u^T = \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}$  :

**Example** The inverse of  $M = I - \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$  is  $M^{-1} = I + \frac{1}{1 - 3} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$

But we don't always start from the identity matrix. Many applications need to invert  $M = A - uv^T$ . After we solve  $Ax = b$  we expect a rank one change to give  $My = b$ . The division by  $1 - v^T u$  above will become a division by  $c = 1 - v^T A^{-1} u = 1 - v^T z$ .

**Step 1** Solve  $Az = u$  and compute  $c = 1 - v^T z$ .

**Step 2** If  $c \neq 0$  then  $M^{-1}b$  is  $y = x + \frac{v^T x}{c} z$ .

Suppose  $A$  is easy to work with.  $A$  might already be factored into  $LU$  by elimination. Then this Sherman-Woodbury-Morrison formula is the fast way to solve  $My = b$ . Here are three problems to end the book !

**9** Take Steps 1–2 to find  $y$  when  $A = I$  and  $u^T = v^T = [1 \ 2 \ 3]$  and  $b^T = [2 \ 1 \ 4]$ .

**10** Step 2 in this “update formula” claims that  $My = (A - uv^T) \left( x + \frac{v^T x}{c} z \right) = b$ .

Simplify this to  $\frac{uv^T x}{c} [1 - c - v^T z] = 0$ . This is true since  $c = 1 - v^T z$ .

**11** When  $A$  has a new row  $v^T$ ,  $A^T A$  in the least squares equation changes to  $M$  :

$$M = \begin{bmatrix} A^T & v \end{bmatrix} \begin{bmatrix} A \\ v^T \end{bmatrix} = A^T A + vv^T = \text{rank one change in } A^T A.$$

Why is that multiplication correct? The updated  $\hat{x}_{\text{new}}$  comes from Steps 1 and 2.

For reference here are four formulas for  $M^{-1}$ . The first two were given above, when the change was  $uv^T$ . Formulas 3 and 4 go beyond rank one to allow matrices  $U, V, W$ .

**1**  $M = I - uv^T$  and  $M^{-1} = I + uv^T / (1 - v^T u)$  (rank1 change)

**2**  $M = A - uv^T$  and  $M^{-1} = A^{-1} + A^{-1} uv^T A^{-1} / (1 - v^T A^{-1} u)$

**3**  $M = I - UV$  and  $M^{-1} = I_n + U(I_m - VU)^{-1} V$

**4**  $M = A - UW^{-1}V$  and  $M^{-1} = A^{-1} + A^{-1} U(W - VA^{-1}U)^{-1} VA^{-1}$

Formula 4 is the “matrix inversion lemma” in engineering. Not seen until now! The Kalman filter for solving block tridiagonal systems uses formula 4 at each step.