# Lecture 6: Discrete Random Variables

#### 19 September 2005

## 1 Expectation

The expectation of a random variable is its average value, with weights in the average given by the probability distribution

$$\mathbf{E}\left[X\right] = \sum_{x} \Pr\left(X = x\right) x$$

If c is a constant,  $\mathbf{E}[c] = c$ . If a and b are constants,  $\mathbf{E}[aX + b] = a\mathbf{E}[X] + b$ . If  $X \ge Y$ , then  $\mathbf{E}[X] \ge \mathbf{E}[Y]$ Now let's think about  $\mathbf{E}[X + Y]$ .

$$\begin{split} \mathbf{E} \left[ X+Y \right] &= \sum_{x,y} \left( x+y \right) \mathrm{Pr} \left( X=x,Y=y \right) \\ &= \sum_{x,y} x \mathrm{Pr} \left( X=x,Y=y \right) + \sum_{x,y} y \mathrm{Pr} \left( X=x,Y=y \right) \\ &= \sum_{x} x \sum_{y} \mathrm{Pr} \left( X=x,Y=y \right) + \sum_{y} y \sum_{x} \mathrm{Pr} \left( X=x,Y=y \right) \end{split}$$

by total probability,  $\sum_x \Pr(X = x, Y = y) = \Pr(X = x)$ , likewise  $\sum_x \Pr(X = x, Y = y) = \Pr(Y = y)$ . So,

$$\mathbf{E} [X + Y] = \sum_{x} x \Pr (X = x) + \sum_{y} y \Pr (Y = y)$$
$$= \mathbf{E} [X] + \mathbf{E} [Y]$$

Notice that  $\mathbf{E}[X]$  works just like a mean; in fact we can think of it as being the population mean (as opposed to the sample mean).

The variance is the expectation of  $(X - \mathbf{E}[X])^2$ .

$$\operatorname{Var}(X) = \sum_{x} p(x)(x - \mathbf{E}[X])^{2}$$

which we can show is  $\mathbf{E}[X^2] - (\mathbf{E}[X])^2$ .

$$\operatorname{Var}(X) = \mathbf{E}\left[\left(X - \mathbf{E}\left[X\right]\right)^{2}\right]$$
$$= \mathbf{E}\left[X^{2} - 2X\mathbf{E}\left[X\right] + \left(\mathbf{E}\left[X\right]\right)^{2}\right]$$
$$= \mathbf{E}\left[X^{2}\right] - \mathbf{E}\left[2X\mathbf{E}\left[X\right]\right] + \mathbf{E}\left[\left(\mathbf{E}\left[X\right]\right)^{2}\right]$$

Now  $\mathbf{E}[X]$  is just another constant, so  $\mathbf{E}\left[(\mathbf{E}[X])^2\right] = (\mathbf{E}[X])^2$ , and  $\mathbf{E}[2X\mathbf{E}[X]] = 2\mathbf{E}[X]\mathbf{E}[X] = 2(\mathbf{E}[X])^2$ . So

$$\operatorname{Var}(X) = \mathbf{E} [X^2] - 2(\mathbf{E} [X])^2 + (\mathbf{E} [X])^2$$
$$= \mathbf{E} [X^2] - (\mathbf{E} [X])^2$$

as promised.

The main rule for variance is this:

$$\operatorname{Var}\left(aX+b\right) = a^{2}\operatorname{Var}\left(X\right)$$

It's not generally true that  $\operatorname{Var}(X + Y) = \operatorname{Var}(X) + \operatorname{Var}(Y)$ ; we'll see when it's true later.

#### 1.1 Some useful results

A basic result about expectations is the **Markov inequality**: if X is a non-negative random variable, and a is a positive constant, then

$$\Pr\left(X \ge a\right) \le \frac{\mathbf{E}\left[X\right]}{a}$$

*Proof*: Let  $A = \{X \ge a\}$ . So  $X \ge a1_A$ : either  $1_A = 0$ , in which case  $X \ge 0$ , or else  $1_A = 1$ , but then  $X \ge a$ . So  $\mathbf{E}[X] \ge \mathbf{E}[a1_A] = a\mathbf{E}[1_A] = a\Pr(X \ge a)$ .

The **Chebyshev inequality** is a special case of the Markov inequality, but a very useful one. It's plain that  $(X - \mathbf{E}[X])^2 \ge 0$ , so applying the Markov inequality gives

$$\Pr\left(\left(X - \mathbf{E}\left[X\right]\right)^2 \ge a^2\right) \le \frac{\operatorname{Var}\left(X\right)}{a^2}$$

Taking the square root of the term inside the left-hand side,

$$\Pr\left(|X - \mathbf{E}[X]| \ge a\right) \le \frac{\operatorname{Var}(X)}{a^2}$$

The Chebyshev inequality helps give meaning to the variance: it tells us about how unlikely it is for the random variable to depart very far from its mean.

## 2 Independent r.v.s

We'll say that random variables are independent if their probability distributions factor,  $\Pr(X = x, Y = y) = \Pr(X = x) \Pr(Y = y)$ .

If the variables are independent, then  $\mathbf{E}[XY] = \mathbf{E}[X]\mathbf{E}[Y]$ .

$$\mathbf{E} [XY] = \sum_{x,y} xy \operatorname{Pr} (X = x, Y = y)$$

$$= \sum_{x} \sum_{y} xy \operatorname{Pr} (X = x, Y = y)$$

$$= \sum_{x} \sum_{y} xy \operatorname{Pr} (X = x) \operatorname{Pr} (Y = y)$$

$$= \sum_{x} x \operatorname{Pr} (X = x) \sum_{y} y \operatorname{Pr} (Y = y)$$

$$= \sum_{x} x \operatorname{Pr} (X = x) \mathbf{E} [Y]$$

$$= \mathbf{E} [Y] \sum_{x} x \operatorname{Pr} (X = x)$$

$$= \mathbf{E} [Y] \mathbf{E} [X]$$

This isn't the only time that  $\mathbf{E}[XY] = \mathbf{E}[X]\mathbf{E}[Y]$ , though.

Here's where independence gets important: what's the variance of X + Y?

$$Var (X + Y) = \mathbf{E} \left[ (X + Y)^2 \right] - (\mathbf{E} [X + Y])^2$$
  
=  $\mathbf{E} \left[ X^2 + 2XY + Y^2 \right] - (\mathbf{E} [X] + \mathbf{E} [Y])^2$   
=  $\mathbf{E} \left[ X^2 \right] + 2\mathbf{E} [XY] + \mathbf{E} \left[ Y^2 \right] - \left[ (\mathbf{E} [X])^2 + 2\mathbf{E} [X] \mathbf{E} [Y] + (\mathbf{E} [Y])^2 \right]$   
=  $\mathbf{E} \left[ X^2 \right] - (\mathbf{E} [X])^2 + \mathbf{E} \left[ Y^2 \right] - (\mathbf{E} [Y])^2 + 2\mathbf{E} [XY] - 2\mathbf{E} [X] \mathbf{E} [Y]$   
=  $Var (X) + Var (Y) + 2\mathbf{E} [XY] - 2\mathbf{E} [X] \mathbf{E} [Y]$ 

But we've just seen that  $\mathbf{E}[XY] = \mathbf{E}[X] \mathbf{E}[Y]$  if X and Y are independent, so then

$$\operatorname{Var}(X+Y) = \operatorname{Var}(X) + \operatorname{Var}(Y)$$

### **3** Binomial random variables

Recall that the distribution of the binomial is

$$ProbX = x = \binom{n}{x} p^x (1-p)^{n-x}$$

and that it's the sum of n independent Bernoulli variables with parameter p.

How do we know this is a valid probability distribution? It's clearly  $\geq 0$  for all x, but how do I know it sums to 1? Because of the binomial theorem from algebra (which is where the name comes from).

$$(a+b)^{n} = \sum_{k=0}^{n} \binom{n}{k} a^{k} b^{n-k}$$
$$(p+(1-p))^{n} = \sum_{k=0}^{n} \binom{n}{k} p^{k} (1-p)^{n-k}$$
$$1^{n} = \sum_{k=0}^{n} \binom{n}{k} p^{k} (1-p)^{n-k}$$
$$1 = \sum_{k=0}^{n} \binom{n}{k} p^{k} (1-p)^{n-k}$$

To find the mean and variance, we could either do the appropriate sums explicitly, which means using ugly tricks about the binomial formula; or we could use the fact that X is a sum of n independent Bernoulli variables. Because the Bernoulli variables have expectation p,  $\mathbf{E}[X] = np$ . Because they have variance p(1-p), Var(X) = np(1-p).

#### 4 Geometric random variables

Suppose we keep trying independent Bernoulli variables until we have a success; each has probability of success p. Then the probability that the number of failures is k is  $(1-p)^k p$ . (Be careful, some people use p as the probability of failure here, i.e. they reverse p and 1-p.)

First, check that this weird thing is a valid probability distribution — does it sum to one? Yes:

$$\sum_{k=0}^{\infty} (1-p)^k p = p \sum_{k=0}^{\infty} (1-p)^k = p \frac{1}{1-(1-p)} = p \frac{1}{p} = 1$$

This uses the geometric series, the fact that  $\sum p^k = 1/(1-p)$ , if p is between 0 and 1.

Now let's think about the mean.

$$\mathbf{E}[X] = \sum_{k=0}^{\infty} k(1-p)^k p = p \sum_{k=1}^{\infty} k(1-p)^k = p \frac{1}{p^2} = \frac{1}{p}$$

Similarly (but more involvedly) the variance is  $(1-p)/p^2$ .

#### 4.1 Negative binomial random variables

Instead of just getting one success, we might keep going until we get r of them. The probability distribution then is just  $\Pr(X = k) = {\binom{k-1}{r-1}}p^r(1-p)^{k-r}, k \ge r$ . If we think of  $W_1$  as the number of trials we have to make to get the first success, and then  $W_2$  the number of further trials to the second success, and so on, we can see that  $X = W_1 + W_2 + \ldots + W_r$ , and that the  $W_i$  are independent and geometric random variables. So  $\mathbf{E}[X] = r/p$ , and  $\operatorname{Var}(X) = r(1-p)/p^2$ .

## 5 Poisson random variables

Think about a very large number of Bernoulli trials, where  $n \to \infty$ , but the expected number of successes stays constant, say  $\lambda$ . For instance, suppose we're looking at the number of particles emitted by a chunk of radioactive substance over one-second intervals of time. Every atom has a certain probability to decay over a given unit of time; as we make the time intervals smaller, we make those probabilities smaller, but the average total should still come to the same number.

If we have only a finite n, but n is very large, so  $p = \lambda/n$ .

$$\Pr(X = k) = \binom{n}{k} p^k (1-p)^k$$
$$= \frac{n!}{k!(n-k)!} \frac{p^k}{(1-p)^k} (1-p)^n$$

Since *n* is large, we can use Stirling's approximation on *n*! and (n - k)!, so  $n! \approx n^n$  and  $(n - k)! \approx (n - k)^{n-k} \approx n^{n-k}$ .

$$\begin{split} \Pr\left(X=k\right) &\approx \quad \frac{n^k}{k!} \frac{p^k}{\left(1-p\right)^k} \left(1-\frac{\lambda}{n}\right)^n \\ &\to \quad \frac{\lambda^k}{k!} e^{-\lambda} \end{split}$$

because  $\lim (1 + x/n)^n = e^x$ .

We can check that the probability adds up to one, because

$$\sum_{k=0}^{\infty} \frac{\lambda^k}{k!} = e^{\lambda}$$

We can also get the mean:

$$\mathbf{E}[X] = \sum_{k=0}^{\infty} k \frac{\lambda^k}{k!} e^{-\lambda}$$
$$= \sum_{k=1}^{\infty} k \frac{\lambda^k}{k!} e^{-\lambda}$$
$$= \sum_{k=1}^{\infty} \lambda \frac{\lambda^{k-1}}{(k-1)!} e^{-\lambda}$$

$$= \lambda \sum_{k=1}^{\infty} \lambda \frac{\lambda^{k-1}}{(k-1)!} e^{-\lambda}$$
$$= \lambda \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} e^{-\lambda}$$
$$= \lambda$$

The easiest way to get the variance is to first calculate  $\mathbf{E}[X(X-1)]$ , because this will let us use the same sort of trick about factorials and the exponential function again.

$$\mathbf{E} [X(X-1)] = \sum_{k=0}^{\infty} k(k-1) \frac{\lambda^k}{k!} e^{-\lambda}$$
$$\mathbf{E} [X^2 - X] = \sum_{k=2}^{\infty} \frac{\lambda^k}{(k-2)!} e^{-\lambda}$$
$$\mathbf{E} [X^2] - \mathbf{E} [X] = \lambda^2 \sum_{k=2}^{\infty} \frac{\lambda^{k-2}}{(k-2)!} e^{-\lambda}$$
$$= \lambda^2 \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} e^{-\lambda}$$
$$= \lambda^2$$

So  $\mathbf{E}[X^2] = \mathbf{E}[X] + \lambda^2$ .

$$\operatorname{Var}(X) = \mathbf{E} [X^{2}] - (\mathbf{E} [X])^{2}$$
$$= \mathbf{E} [X] + \lambda^{2} - \lambda^{2}$$
$$= \mathbf{E} [X] = \lambda$$

#### Adding Many Independent Random Variables 6

Remember the Chebyshev inequality:

$$\Pr\left(\left|X - \mathbf{E}\left[X\right]\right| \ge a\right) \le \frac{\operatorname{Var}\left(X\right)}{a^2}$$

Let's look at the sum of a whole bunch of independent random variables with the same distribution,  $S_n = \sum_{i=1}^n X_i$ . We know that  $\mathbf{E}[S_n] = \mathbf{E}[\sum_{i=1}^n X_i] = \sum \mathbf{E}[X_i] = n\mathbf{E}[X_1]$ , because they all have the same expectation. Because they're independent, and all have the same variance,  $\operatorname{Var}(S_n) = n\operatorname{Var}(X_1)$ . So

$$\Pr\left(\left|S_{n} - n\mathbf{E}\left[X_{1}\right]\right| \ge a\right) \le \frac{n\operatorname{Var}\left(X_{1}\right)}{a^{2}}$$

Now, notice that  $S_n/n$  is just the sample mean, if we draw a sample of size n. So we can use the Chebyshev inequality to estimate the chance that the sample mean is far from the true, population mean, which is the expectation.

$$\Pr\left(\left|\frac{S_n}{n} - \mathbf{E}[X_1]\right| \ge \epsilon\right) = \Pr\left(|S_n - n\mathbf{E}[X_1]| \ge n\epsilon\right)$$
$$\le \frac{n\operatorname{Var}(X_1)}{n^2\epsilon^2}$$
$$\le \frac{\operatorname{Var}(X_1)}{n\epsilon^2}$$

Observe that whatever  $\epsilon$  is, the probability must go to zero like 1/n (or faster). So the probability that the sample mean differs from the population mean by as much as  $\epsilon$  can be made arbitrarily small, by taking a large enough sample. This is the **law of large numbers**.