Mathematics for Informatics 4a

José Figueroa-O'Farrill



Lecture 9 15 February 2012

→ Ξ

 Discrete random variables X₁,..., X_n on the same probability space have a joint probability mass function:

 $f_{X_1,\ldots,X_n}(x_1,\ldots,x_n)=\mathbb{P}(\{X_1=x_1\}\cap\cdots\cap\{X_n=x_n\})$

 Discrete random variables X₁,..., X_n on the same probability space have a joint probability mass function:

 $f_{X_1,\ldots,X_n}(x_1,\ldots,x_n)=\mathbb{P}(\{X_1=x_1\}\cap\cdots\cap\{X_n=x_n\})$

•
$$f_{X_1,...,X_n} : \mathbb{R}^n \to [0,1]$$
 and $\sum_{x_1,...,x_n} f_{X_1,...,X_n}(x_1,...,x_n) = 1$

 Discrete random variables X₁,..., X_n on the same probability space have a joint probability mass function:

 $f_{X_1,...,X_n}(x_1,\ldots,x_n) = \mathbb{P}(\{X_1 = x_1\} \cap \cdots \cap \{X_n = x_n\})$

• $f_{X_1,...,X_n} : \mathbb{R}^n \to [0, 1]$ and $\sum_{x_1,...,x_n} f_{X_1,...,X_n}(x_1, \dots, x_n) = 1$ • X_1, \dots, X_n are **independent** if for all $2 \leq k \leq n$ and x_{i_1}, \dots, x_{i_k} ,

$$f_{X_{\mathfrak{i}_1},\ldots,X_{\mathfrak{i}_k}}(x_{\mathfrak{i}_1},\ldots,x_{\mathfrak{i}_k}) = f_{X_{\mathfrak{i}_1}}(x_{\mathfrak{i}_1})\ldots f_{X_{\mathfrak{i}_k}}(x_{\mathfrak{i}_k})$$

 Discrete random variables X₁,..., X_n on the same probability space have a joint probability mass function:

 $f_{X_1,...,X_n}(x_1,\ldots,x_n) = \mathbb{P}(\{X_1 = x_1\} \cap \cdots \cap \{X_n = x_n\})$

• $f_{X_1,...,X_n} : \mathbb{R}^n \to [0, 1]$ and $\sum_{x_1,...,x_n} f_{X_1,...,X_n}(x_1,...,x_n) = 1$ • $X_1,...,X_n$ are **independent** if for all $2 \leq k \leq n$ and $x_{i_1},...,x_{i_k}$,

$$f_{X_{i_1},\ldots,X_{i_k}}(x_{i_1},\ldots,x_{i_k}) = f_{X_{i_1}}(x_{i_1})\ldots f_{X_{i_k}}(x_{i_k})$$

• $h(X_1, ..., X_n)$ is a discrete random variable and $\mathbb{E}(h(X_1, ..., X_n)) = \sum_{x_1, ..., x_n} h(x_1, ..., x_n) f_{X_1, ..., X_n}(x_1, ..., x_n)$

 Discrete random variables X₁,..., X_n on the same probability space have a joint probability mass function:

 $f_{X_1,\ldots,X_n}(x_1,\ldots,x_n)=\mathbb{P}(\{X_1=x_1\}\cap\cdots\cap\{X_n=x_n\})$

• $f_{X_1,...,X_n} : \mathbb{R}^n \to [0, 1]$ and $\sum_{x_1,...,x_n} f_{X_1,...,X_n}(x_1,...,x_n) = 1$ • $X_1,...,X_n$ are **independent** if for all $2 \leq k \leq n$ and $x_{i_1},...,x_{i_k}$,

$$f_{X_{i_1},\ldots,X_{i_k}}(x_{i_1},\ldots,x_{i_k}) = f_{X_{i_1}}(x_{i_1})\ldots f_{X_{i_k}}(x_{i_k})$$

- $h(X_1, ..., X_n)$ is a discrete random variable and $\mathbb{E}(h(X_1, ..., X_n)) = \sum_{x_1, ..., x_n} h(x_1, ..., x_n) f_{X_1, ..., X_n}(x_1, ..., x_n)$
- Expectation is linear: $\mathbb{E}(\sum_{i} \alpha_{i} X_{i}) = \sum_{i} \alpha_{i} \mathbb{E}(X_{i})$

Expectation of a product

Lemma

If X and Y are independent, $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$.

< □ > < □ > < □ > < Ξ > < Ξ > < Ξ = ・ ○ < ○

Expectation of a product

Lemma

If X and Y are independent, $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$.

Proof.

$$\mathbb{E}(XY) = \sum_{x,y} xyf_{X,Y}(x,y)$$

= $\sum_{x,y} xyf_X(x)f_Y(y)$ (independence)
= $\sum_x xf_X(x) \sum_y yf_Y(y)$
= $\mathbb{E}(X)\mathbb{E}(Y)$

▲□ ▶ ▲ ■ ▶ ▲ ■ ▶ ■ ■ ■ ● ● ●

The expectation value defines a real inner product.

The expectation value defines a real inner product. If X, Y are two discrete random variables, let us define $\langle X, Y \rangle$ by

 $\langle \mathbf{X}, \mathbf{Y} \rangle = \mathbb{E}(\mathbf{X}\mathbf{Y})$

▲□ ▶ ▲ ■ ▶ ▲ ■ ▶ ■ ■ ■ ● ● ●

The expectation value defines a real inner product. If X, Y are two discrete random variables, let us define $\langle X, Y \rangle$ by

 $\langle \mathbf{X}, \mathbf{Y} \rangle = \mathbb{E}(\mathbf{X}\mathbf{Y})$

We need to show that $\langle X, Y \rangle$ satisfies the axioms of an inner product:

The expectation value defines a real inner product. If X, Y are two discrete random variables, let us define $\langle X, Y \rangle$ by

 $\langle \mathbf{X}, \mathbf{Y} \rangle = \mathbb{E}(\mathbf{X}\mathbf{Y})$

We need to show that $\langle X, Y \rangle$ satisfies the axioms of an inner product:

1 it is symmetric: $\langle X, Y \rangle = \mathbb{E}(XY) = \mathbb{E}(YX) = \langle Y, X \rangle$

4/23

The expectation value defines a real inner product. If X, Y are two discrete random variables, let us define $\langle X, Y \rangle$ by

 $\langle \mathbf{X}, \mathbf{Y} \rangle = \mathbb{E}(\mathbf{X}\mathbf{Y})$

We need to show that $\langle X, Y \rangle$ satisfies the axioms of an inner product:

- **1** it is symmetric: $\langle X, Y \rangle = \mathbb{E}(XY) = \mathbb{E}(YX) = \langle Y, X \rangle$
- it is bilinear:

The expectation value defines a real inner product. If X, Y are two discrete random variables, let us define $\langle X, Y \rangle$ by

 $\langle \mathbf{X}, \mathbf{Y} \rangle = \mathbb{E}(\mathbf{X}\mathbf{Y})$

We need to show that $\langle X, Y \rangle$ satisfies the axioms of an inner product:

- **1** it is symmetric: $\langle X, Y \rangle = \mathbb{E}(XY) = \mathbb{E}(YX) = \langle Y, X \rangle$
- it is bilinear:
 - $\langle aX, Y \rangle = \mathbb{E}(aXY) = a\mathbb{E}(XY) = a \langle X, Y \rangle$

The expectation value defines a real inner product. If X, Y are two discrete random variables, let us define $\langle X, Y \rangle$ by

 $\langle \mathbf{X}, \mathbf{Y} \rangle = \mathbb{E}(\mathbf{X}\mathbf{Y})$

We need to show that $\langle X, Y \rangle$ satisfies the axioms of an inner product:

- **1** it is symmetric: $\langle X, Y \rangle = \mathbb{E}(XY) = \mathbb{E}(YX) = \langle Y, X \rangle$
- it is bilinear:
 - $\langle aX, Y \rangle = \mathbb{E}(aXY) = a\mathbb{E}(XY) = a \langle X, Y \rangle$
 - $\langle X_1 + X_2, Y \rangle = \mathbb{E}((X_1 + X_2)Y) = \mathbb{E}(X_1Y) + \mathbb{E}(X_2Y) = \langle X_1, Y \rangle + \langle X_2, Y \rangle$

The expectation value defines a real inner product. If X, Y are two discrete random variables, let us define $\langle X, Y \rangle$ by

 $\langle \mathbf{X}, \mathbf{Y} \rangle = \mathbb{E}(\mathbf{X}\mathbf{Y})$

We need to show that $\langle X,Y\rangle$ satisfies the axioms of an inner product:

- **1** it is symmetric: $\langle X, Y \rangle = \mathbb{E}(XY) = \mathbb{E}(YX) = \langle Y, X \rangle$
- it is bilinear:
 - $\langle aX, Y \rangle = \mathbb{E}(aXY) = a\mathbb{E}(XY) = a \langle X, Y \rangle$
 - $\langle X_1 + X_2, Y \rangle = \mathbb{E}((X_1 + X_2)Y) = \mathbb{E}(X_1Y) + \mathbb{E}(X_2Y) = \langle X_1, Y \rangle + \langle X_2, Y \rangle$

3 it is positive-definite: if $\langle X, X \rangle = 0$, then $\mathbb{E}(X^2) = 0$, whence $\sum_x x^2 f(x) = 0$, whence x f(x) = 0 for all x. If $x \neq 0$, then f(x) = 0 and thus f(0) = 1. In other words, $\mathbb{P}(X = 0) = 1$ and hence X = 0 almost surely.

Additivity of variance for independent variables

How about the variance Var(X + Y)?

Additivity of variance for independent variables

How about the variance Var(X + Y)?

 $Var(X + Y) = \mathbb{E}((X + Y)^2) - \mathbb{E}(X + Y)^2$ = $\mathbb{E}(X^2 + 2XY + Y^2) - (\mathbb{E}(X) + \mathbb{E}(Y))^2$ = $\mathbb{E}(X^2) + 2\mathbb{E}(XY) + \mathbb{E}(Y^2) - \mathbb{E}(X)^2 - 2\mathbb{E}(X)\mathbb{E}(Y) - \mathbb{E}(Y)^2$ = $Var(X) + Var(Y) + 2(\mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y))$

Additivity of variance for independent variables

How about the variance Var(X + Y)?

$$Var(X + Y) = \mathbb{E}((X + Y)^2) - \mathbb{E}(X + Y)^2$$

= $\mathbb{E}(X^2 + 2XY + Y^2) - (\mathbb{E}(X) + \mathbb{E}(Y))^2$
= $\mathbb{E}(X^2) + 2\mathbb{E}(XY) + \mathbb{E}(Y^2) - \mathbb{E}(X)^2 - 2\mathbb{E}(X)\mathbb{E}(Y) - \mathbb{E}(Y)^2$
= $Var(X) + Var(Y) + 2(\mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y))$

Theorem

If X and Y are independent discrete random variables

Var(X + Y) = Var(X) + Var(Y)

Covariance

Definition

The **covariance** of two discrete random variables is

 $\mathsf{Cov}(X, Y) = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y)$

< ロ > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 >

Covariance

Definition

The **covariance** of two discrete random variables is

 $\operatorname{Cov}(X, Y) = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y)$

Letting μ_X and μ_Y denote the means of X and Y, respectively,

 $\text{Cov}(X,Y) = \mathbb{E}((X-\mu_X)(Y-\mu_Y))$

Covariance

Definition

The covariance of two discrete random variables is

 $\operatorname{Cov}(X, Y) = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y)$

Letting μ_X and μ_Y denote the means of X and Y, respectively,

$$\operatorname{Cov}(X, Y) = \mathbb{E}((X - \mu_X)(Y - \mu_Y))$$

Indeed,

$$\begin{split} \mathbb{E}((X-\mu_X)(Y-\mu_Y)) &= \mathbb{E}(XY) - \mathbb{E}(\mu_XY) - \mathbb{E}(\mu_YX) + \mathbb{E}(\mu_X\mu_Y) \\ &= \mathbb{E}(XY) - \mu_X\mu_Y \end{split}$$

We roll two fair dice. Let X and Y denote their scores. Independence says that Cov(X, Y) = 0.

▲□ ▶ ▲ ■ ▶ ▲ ■ ▶ ■ ■ ■ ● ● ●

We roll two fair dice. Let X and Y denote their scores. Independence says that Cov(X, Y) = 0. Consider however the new variables U = min(X, Y) and V = max(X, Y):

We roll two fair dice. Let X and Y denote their scores. Independence says that Cov(X, Y) = 0. Consider however the new variables U = min(X, Y) and V = max(X, Y):

u	1	2	3	4	5	6
1	1	1	1	1	1	1
2	1	2	2	2	2	2
3	1	2	3	3	3	3
4	1	2	3	4	4	4
5	1	2	3	4	5	5
6	1	2	3	4	5	6

▲□ ▶ ▲ ■ ▶ ▲ ■ ▶ ■ ■ ■ ● ● ●

José Figueroa-O'Farrill

We roll two fair dice. Let X and Y denote their scores. Independence says that Cov(X, Y) = 0. Consider however the new variables U = min(X, Y) and V = max(X, Y):

u	1	2	3	4	5	6	V	1	2	3	4	5	6
1	1	1	1	1	1	1	1	1	2	3	4	5	6
2	1	2	2	2	2	2	2	2	2	3	4	5	6
3	1	2	3	3	3	3	3	3	3	3	4	5	6
4	1	2	3	4	4	4	4	4	4	4	4	5	6
5	1	2	3	4	5	5	5	5	5	5	5	5	6
6	1	2	3	4	5	6	6	6	6	6	6	6	6

(日本)

We roll two fair dice. Let X and Y denote their scores. Independence says that Cov(X, Y) = 0. Consider however the new variables U = min(X, Y) and V = max(X, Y):

u	1	2	3	4	5	6		V	1	2	3	4	5	6
1	1	1	1	1	1	1		1	1	2	3	4	5	6
2	1	2	2	2	2	2		2	2	2	3	4	5	6
3	1	2	3	3	3	3		3	3	3	3	4	5	6
4	1	2	3	4	4	4		4	4	4	4	4	5	6
5	1	2	3	4	5	5		5	5	5	5	5	5	6
6	1	2	3	4	5	6		6	6	6	6	6	6	6
$\mathbb{E}(\mathbf{U}) = \frac{91}{36}, \ \mathbb{E}(\mathbf{U}^2) = \frac{301}{36}, \ \mathbb{E}(\mathbf{V}) = \frac{161}{36}, \ \mathbb{E}(\mathbf{V}^2) = \frac{791}{36}, \ \mathbb{E}(\mathbf{U}\mathbf{V}) = \frac{49}{4}$														

7/23

We roll two fair dice. Let X and Y denote their scores. Independence says that Cov(X, Y) = 0. Consider however the new variables U = min(X, Y) and V = max(X, Y):

u	1	2	3	4	5	6		V	1	2	3	4	5	6
1	1	1	1	1	1	1	-	1	1	2	3	4	5	6
2	1	2	2	2	2	2		2	2	2	3	4	5	6
3	1	2	3	3	3	3		3	3	3	3	4	5	6
4	1	2	3	4	4	4		4	4	4	4	4	5	6
5	1	2	3	4	5	5		5	5	5	5	5	5	6
6	1	2	3	4	5	6		6	6	6	6	6	6	6
$\mathbb{E}(\mathbf{U}) = \frac{91}{36}, \ \mathbb{E}(\mathbf{U}^2) = \frac{301}{36}, \ \mathbb{E}(\mathbf{V}) = \frac{161}{36}, \ \mathbb{E}(\mathbf{V}^2) = \frac{791}{36}, \ \mathbb{E}(\mathbf{U}\mathbf{V}) = \frac{49}{4}$ $\implies Var(\mathbf{U}) = Var(\mathbf{V}) = \frac{2555}{1296} \text{and} Cov(\mathbf{U}, \mathbf{V}) = \left(\frac{35}{36}\right)^2$														

▲□ ▶ ▲ □ ▶ ▲ □ ▶ □ □ ● ● ●

Two discrete random variables X and Y are said to be **uncorrelated** if Cov(X, Y) = 0.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三回日 のへで

Two discrete random variables X and Y are said to be **uncorrelated** if Cov(X, Y) = 0.

Warning

Uncorrelated random variables need not be independent!

Two discrete random variables X and Y are said to be **uncorrelated** if Cov(X, Y) = 0.

Warning

Uncorrelated random variables need not be independent!

Counterexample

Suppose that X is a discrete random variable with probability mass function symmetric about 0; that is, $f_X(-x) = f_X(x)$. Let $Y = X^2$.

Two discrete random variables X and Y are said to be **uncorrelated** if Cov(X, Y) = 0.

Warning

Uncorrelated random variables need not be independent!

Counterexample

Suppose that X is a discrete random variable with probability mass function symmetric about 0; that is, $f_X(-x) = f_X(x)$. Let $Y = X^2$. Clearly X, Y are not independent: f(x, y) = 0 unless $y = x^2$ even if $f_X(x)f_Y(y) \neq 0$.

Two discrete random variables X and Y are said to be **uncorrelated** if Cov(X, Y) = 0.

Warning

Uncorrelated random variables need not be independent!

Counterexample

Suppose that X is a discrete random variable with probability mass function symmetric about 0; that is, $f_X(-x) = f_X(x)$. Let $Y = X^2$. Clearly X, Y are not independent: f(x, y) = 0 unless $y = x^2$ even if $f_X(x)f_Y(y) \neq 0$. However they are uncorrelated:

$$\mathbb{E}(XY) = \mathbb{E}(X^3) = \sum_{x} x^3 f_X(x) = 0$$

Two discrete random variables X and Y are said to be **uncorrelated** if Cov(X, Y) = 0.

Warning

Uncorrelated random variables need not be independent!

Counterexample

Suppose that X is a discrete random variable with probability mass function symmetric about 0; that is, $f_X(-x) = f_X(x)$. Let $Y = X^2$. Clearly X, Y are not independent: f(x, y) = 0 unless $y = x^2$ even if $f_X(x)f_Y(y) \neq 0$. However they are uncorrelated:

$$\mathbb{E}(XY) = \mathbb{E}(X^3) = \sum_{x} x^3 f_X(x) = \mathbf{0}$$

and similarly $\mathbb{E}(X) = 0$, whence $\mathbb{E}(X)\mathbb{E}(Y) = 0$.

An alternative criterion for independence

The above counterexample says that the following implication cannot be reversed:

X, Y independent $\implies \mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$

▲■▶ ▲■▶ ▲■▶ 差目目 のへで

An alternative criterion for independence

The above counterexample says that the following implication cannot be reversed:

X, Y independent $\implies \mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$

However, one has the following

Theorem

Two discrete random variables \boldsymbol{X} and \boldsymbol{Y} are independent if and only if

 $\mathbb{E}(g(X)h(Y)) = \mathbb{E}(g(X))\mathbb{E}(h(Y))$

for all functions g, h.

▲□ ▶ ▲ □ ▶ ▲ □ ▶ □ □ ● ● ●

9/23
An alternative criterion for independence

The above counterexample says that the following implication cannot be reversed:

X, Y independent $\implies \mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$

However, one has the following

Theorem

Two discrete random variables \boldsymbol{X} and \boldsymbol{Y} are independent if and only if

 $\mathbb{E}(g(X)h(Y)) = \mathbb{E}(g(X))\mathbb{E}(h(Y))$

for all functions g, h.

The proof is not hard, but we will skip it.

▲□ ▶ ▲ □ ▶ ▲ □ ▶ □ □ ● ● ●

Recall that $\langle X, Y \rangle = \mathbb{E}(XY)$ is an inner product.

▲□ ▶ ▲ □ ▶ ▲ □ ▶ □ □ ● ● ● ●

10/23

Recall that $\langle X, Y \rangle = \mathbb{E}(XY)$ is an inner product. Every inner product obeys the \checkmark Cauchy–Schwarz inequality:

 $\left< X,Y \right>^2 \leqslant \left< X,X \right> \left< Y,Y \right>$

通 ト イヨ ト イヨ ト 三日日 のくべ

Recall that $\langle X, Y \rangle = \mathbb{E}(XY)$ is an inner product. Every inner product obeys the Cauchy-Schwarz inequality:

 $\left< X,Y \right>^2 \leqslant \left< X,X \right> \left< Y,Y \right>$

which in terms of expectations is

 $\mathbb{E}(XY)^2 \leqslant \mathbb{E}(X^2)\mathbb{E}(Y^2)$

▲□ ▶ ▲ □ ▶ ▲ □ ▶ □ □ ● ● ●

Recall that $\langle X, Y \rangle = \mathbb{E}(XY)$ is an inner product. Every inner product obeys the Cauchy-Schwarz inequality:

 $\left< X,Y \right>^2 \leqslant \left< X,X \right> \left< Y,Y \right>$

which in terms of expectations is

 $\mathbb{E}(XY)^2 \leqslant \mathbb{E}(X^2)\mathbb{E}(Y^2)$

Now,

 $\text{Cov}(X,Y)^2 = \mathbb{E}((X-\mu_X)(Y-\mu_Y))^2 \leqslant \mathbb{E}((X-\mu_X)^2)\mathbb{E}((Y-\mu_Y)^2)$

Recall that $\langle X, Y \rangle = \mathbb{E}(XY)$ is an inner product. Every inner product obeys the Cauchy-Schwarz inequality:

 $\left< X,Y \right>^2 \leqslant \left< X,X \right> \left< Y,Y \right>$

which in terms of expectations is

 $\mathbb{E}(XY)^2 \leqslant \mathbb{E}(X^2)\mathbb{E}(Y^2)$

Now,

 $\text{Cov}(X,Y)^2 = \mathbb{E}((X-\mu_X)(Y-\mu_Y))^2 \leqslant \mathbb{E}((X-\mu_X)^2)\mathbb{E}((Y-\mu_Y)^2)$

whence

$$Cov(X, Y)^2 \leq Var(X) Var(Y)$$

Let X and Y be two discrete random variables with means μ_X and μ_Y and standard deviations σ_X , σ_Y .

Let X and Y be two discrete random variables with means μ_X and μ_Y and standard deviations σ_X , σ_Y . The correlation $\rho(X, Y)$ of X and Y is defined by

 $\rho(X,Y) = \frac{\mathsf{Cov}(X,Y)}{\sigma_X \sigma_Y}$

Let X and Y be two discrete random variables with means μ_X and μ_Y and standard deviations σ_X , σ_Y . The **correlation** $\rho(X, Y)$ of X and Y is defined by

 $\rho(X,Y) = \frac{\mathsf{Cov}(X,Y)}{\sigma_X \sigma_Y}$

From the Cauchy-Schwarz inequality, we see that

 $\rho(X,Y)^2 \leqslant 1 \implies -1 \leqslant \rho(X,Y) \leqslant 1$

Let X and Y be two discrete random variables with means μ_X and μ_Y and standard deviations σ_X , σ_Y . The **correlation** $\rho(X, Y)$ of X and Y is defined by

 $\rho(X,Y) = \frac{\mathsf{Cov}(X,Y)}{\sigma_X \sigma_Y}$

From the Cauchy-Schwarz inequality, we see that

 $\rho(X,Y)^2 \leqslant 1 \implies -1 \leqslant \rho(X,Y) \leqslant 1$

Hence the correlation is a number between -1 and 1:

Let X and Y be two discrete random variables with means μ_X and μ_Y and standard deviations σ_X , σ_Y . The **correlation** $\rho(X, Y)$ of X and Y is defined by

 $\rho(X,Y) = \frac{\mathsf{Cov}(X,Y)}{\sigma_X \sigma_Y}$

From the Cauchy-Schwarz inequality, we see that

 $\rho(X,Y)^2 \leqslant 1 \implies -1 \leqslant \rho(X,Y) \leqslant 1$

Hence the correlation is a number between -1 and 1:

 a correlation of 1 suggests a linear relation with positive slope between X and Y,

Let X and Y be two discrete random variables with means μ_X and μ_Y and standard deviations σ_X , σ_Y . The **correlation** $\rho(X, Y)$ of X and Y is defined by

 $\rho(X,Y) = \frac{\mathsf{Cov}(X,Y)}{\sigma_X \sigma_Y}$

From the Cauchy-Schwarz inequality, we see that

 $\rho(X,Y)^2 \leqslant 1 \implies -1 \leqslant \rho(X,Y) \leqslant 1$

Hence the correlation is a number between -1 and 1:

- a correlation of 1 suggests a linear relation with positive slope between X and Y,
- whereas a correlation of -1 suggests a linear relation with negative slope.

Example (Max and min for two fair dice – continued) Continuing with the revious example, we now simply compute $\rho(U, V) = \frac{\text{Cov}(U, V)}{\sqrt{\text{Var}(U) \text{Var}(V)}} = \frac{35^2}{36^2} / \frac{2555}{36^2} = \frac{35}{73}.$

Example (Max and min for two fair dice – continued)

Continuing with the previous example, we now simply compute

$$\rho(\mathbf{U}, \mathbf{V}) = \frac{\text{Cov}(\mathbf{U}, \mathbf{V})}{\sqrt{\text{Var}(\mathbf{U}) \text{Var}(\mathbf{V})}} = \frac{35^2}{36^2} \left/ \frac{2555}{36^2} = \frac{35}{75} \right|_{10}$$

Remark

The funny normalisation of $\rho(X, Y)$ is justified by the following:

 $\rho(\alpha X + \beta, \gamma Y + \delta) = \text{sign}(\alpha \gamma) \rho(X, Y)$

Example (Max and min for two fair dice – continued)

Continuing with the previous example, we now simply compute

$$\rho(\mathbf{U},\mathbf{V}) = \frac{\text{Cov}(\mathbf{U},\mathbf{V})}{\sqrt{\text{Var}(\mathbf{U})\text{ Var}(\mathbf{V})}} = \frac{35^2}{36^2} \left/ \frac{2555}{36^2} = \frac{35}{73}.$$

Remark

The funny normalisation of $\rho(X, Y)$ is justified by the following:

 $\rho(\alpha X + \beta, \gamma Y + \delta) = \text{sign}(\alpha \gamma)\rho(X, Y)$

which follows from

 $Cov(\alpha X + \beta, \gamma Y + \delta) = \alpha \gamma Cov(X, Y)$

Example (Max and min for two fair dice – continued)

Continuing with the previous example, we now simply compute

$$\rho(\mathbf{U},\mathbf{V}) = \frac{\text{Cov}(\mathbf{U},\mathbf{V})}{\sqrt{\text{Var}(\mathbf{U})\text{ Var}(\mathbf{V})}} = \frac{35^2}{36^2} \left/ \frac{2555}{36^2} = \frac{35}{73}.$$

Remark

The funny normalisation of $\rho(X, Y)$ is justified by the following:

 $\rho(\alpha X + \beta, \gamma Y + \delta) = \text{sign}(\alpha \gamma)\rho(X, Y)$

which follows from

 $\operatorname{Cov}(\alpha X + \beta, \gamma Y + \delta) = \alpha \gamma \operatorname{Cov}(X, Y)$

and $\sigma_{\alpha X+\beta} = |\alpha|\sigma_X$ and $\sigma_{\gamma Y+\delta} = |\gamma|\sigma_Y$.

Markov's inequality

Theorem (Markov's inequality)

◎ ▶ ▲ 臣 ▶ ▲ 臣 ▶ 三 臣 ■ の Q @

Markov's inequality

Theorem (Markov's inequality)

Let X be a discrete random variable taking non-negative values. Then

 $\mathbb{P}(X \ge a) \leqslant \frac{\mathbb{E}(X)}{a}$



José Figueroa-O'Farrill mi4a (Probability) Lecture 9

▲□ ▶ ▲ ■ ▶ ▲ ■ ▶ ■ ■ ■ ● ● ●

Markov's inequality

Theorem (Markov's inequality)

Let X be a discrete random variable taking non-negative values. Then

 $\mathbb{P}(X \ge a) \leqslant \frac{\mathbb{E}(X)}{a}$



Proof.

$$\mathbb{E}(X) = \sum_{x \ge 0} x \mathbb{P}(X = x) = \sum_{0 \le x < a} x \mathbb{P}(X = x) + \sum_{x \ge a} x \mathbb{P}(X = x)$$
$$\geq \sum_{x \ge a} x \mathbb{P}(X = x) \ge \sum_{x \ge a} a \mathbb{P}(X = x) = a \mathbb{P}(X \ge a)$$

A factory produces an average of n items every week. What can be said about the probability that this week's production shall be at least 2n items?

A = A = A = E = OQO

A factory produces an average of n items every week. What can be said about the probability that this week's production shall be at least 2n items?

Let X be the discrete random variable counting the number of items produced. Then by Markov's inequality

$$\mathbb{P}(X \geqslant 2n) \leqslant \frac{n}{2n} = \frac{1}{2} .$$

A = A = A = A = A < A
</p>

A factory produces an average of n items every week. What can be said about the probability that this week's production shall be at least 2n items?

Let X be the discrete random variable counting the number of items produced. Then by Markov's inequality

$$\mathbb{P}(X \geqslant 2n) \leqslant \frac{n}{2n} = \frac{1}{2} .$$

So I wouldn't bet on it!

A factory produces an average of n items every week. What can be said about the probability that this week's production shall be at least 2n items?

Let X be the discrete random variable counting the number of items produced. Then by Markov's inequality

$$\mathbb{P}(X \geqslant 2n) \leqslant \frac{n}{2n} = \frac{1}{2} .$$

So I wouldn't bet on it!

Markov's inequality is not terribly sharp; e.g.,

 $\mathbb{P}(X \geqslant \mathbb{E}(X)) \leqslant 1$.

同ト (ヨト (ヨト ヨヨ) の(へ)

A factory produces an average of n items every week. What can be said about the probability that this week's production shall be at least 2n items?

Let X be the discrete random variable counting the number of items produced. Then by Markov's inequality

$$\mathbb{P}(X \geqslant 2n) \leqslant \frac{n}{2n} = \frac{1}{2} .$$

So I wouldn't bet on it!

Markov's inequality is not terribly sharp; e.g.,

 $\mathbb{P}(X \geqslant \mathbb{E}(X)) \leqslant 1$.

It has one interesting corollary, though.





< □ > < □ > < 三 > < 三 > < 三 > < 三 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

José Figueroa-O'Farrill mi4a (Probability) Lecture 9

Theorem

Let X be a discrete random variable with mean μ and variance σ^2 . Then for any $\epsilon > 0$,

$$\mathbb{P}(|\mathsf{X}-\boldsymbol{\mu}| \ge \varepsilon) \leqslant \frac{\sigma^2}{\varepsilon^2}$$



15/23

Theorem

Let X be a discrete random variable with mean μ and variance σ^2 . Then for any $\epsilon > 0$,

 $\mathbb{P}(|\mathbf{X} - \boldsymbol{\mu}| \ge \varepsilon) \le \frac{\sigma^2}{\varepsilon^2}$



Proof.

Notice that for $\varepsilon > 0$, $|X - \mu| \ge \varepsilon$ if and only if $(X - \mu)^2 \ge \varepsilon^2$,

Theorem

Let X be a discrete random variable with mean μ and variance σ^2 . Then for any $\epsilon > 0$,

 $\mathbb{P}(|\mathbf{X} - \boldsymbol{\mu}| \ge \varepsilon) \le \frac{\sigma^2}{\varepsilon^2}$



Proof.

Notice that for $\varepsilon > 0$, $|X - \mu| \ge \varepsilon$ if and only if $(X - \mu)^2 \ge \varepsilon^2$, so

$$\begin{split} \mathbb{P}(|X-\mu| \geqslant \epsilon) &= \mathbb{P}((X-\mu)^2 \geqslant \epsilon^2) \\ &\leqslant \frac{\mathbb{E}((X-\mu)^2)}{\epsilon^2} = \frac{\sigma^2}{\epsilon^2} \end{split} \tag{by Markov's)} \end{split}$$

Back to the factory in the previous example, let the average be n = 500 and the variance in a week's production is 100, then what can be said about the probability that this week's production falls between 400 and 600?

ゆ くち くち くち して くくろ

Back to the factory in the previous example, let the average be n = 500 and the variance in a week's production is 100, then what can be said about the probability that this week's production falls between 400 and 600? By Chebyshev's,

$$\mathbb{P}(|X-500| \ge 100) \leqslant \frac{\sigma^2}{100^2} = \frac{1}{100}$$

同ト (ヨト (ヨト ヨヨ) の(へ)

Back to the factory in the previous example, let the average be n = 500 and the variance in a week's production is 100, then what can be said about the probability that this week's production falls between 400 and 600? By Chebyshev's,

$$\mathbb{P}(|X-500| \ge 100) \leqslant \frac{\sigma^2}{100^2} = \frac{1}{100}$$

whence

$$\begin{split} \mathbb{P}(|X-500| < 100) &= 1 - \mathbb{P}(|X-500| \geqslant 100) \\ &\geqslant 1 - \frac{1}{100} = \frac{99}{100} \; . \end{split}$$

So pretty likely!

<<p>(日本)

Consider a number n of independent discrete random variables X_1, \ldots, X_n with the same probability mass function.

Consider a number n of independent discrete random variables X_1, \ldots, X_n with the same probability mass function. One says that they are "**independent and identically distributed**", abbreviated "**i.i.d.**".

17/23

Consider a number n of independent discrete random variables X_1, \ldots, X_n with the same probability mass function. One says that they are "**independent and identically distributed**", abbreviated "**i.i.d.**". In particular, they have the same mean and variance.

17/23

Consider a number n of independent discrete random variables X_1, \ldots, X_n with the same probability mass function. One says that they are "**independent and identically distributed**", abbreviated "**i.i.d.**". In particular, they have the same mean and variance.

The law of large numbers says that in the limit $n \to \infty$,

 $\frac{1}{n}\left(X_1+\cdots+X_n\right)\to\mu$

in probability.

Consider a number n of independent discrete random variables X_1, \ldots, X_n with the same probability mass function. One says that they are "**independent and identically distributed**", abbreviated "**i.i.d.**". In particular, they have the same mean and variance.

The law of large numbers says that in the limit $n \to \infty$,

 $\frac{1}{n}\left(X_1+\cdots+X_n\right)\to\mu$

in probability.

The law of large numbers justifies the "relative frequency" interpretation of probability.
Consider a number n of independent discrete random variables X_1, \ldots, X_n with the same probability mass function. One says that they are "**independent and identically distributed**", abbreviated "**i.i.d.**". In particular, they have the same mean and variance.

The law of large numbers says that in the limit $n \to \infty$,

 $\frac{1}{n}\left(X_1+\cdots+X_n\right)\to\mu$

in probability.

The law of large numbers justifies the "relative frequency" interpretation of probability. For example, it says that tossing a fair coin a large number n of times, the proportion of heads will approach $\frac{1}{2}$ in the limit $n \to \infty$,

Consider a number n of independent discrete random variables X_1, \ldots, X_n with the same probability mass function. One says that they are "**independent and identically distributed**", abbreviated "**i.i.d.**". In particular, they have the same mean and variance.

The law of large numbers says that in the limit $n \to \infty$,

 $\frac{1}{n}\left(X_1+\cdots+X_n\right)\to\mu$

in probability.

The law of large numbers justifies the "relative frequency" interpretation of probability. For example, it says that tossing a fair coin a large number n of times, the proportion of heads will approach $\frac{1}{2}$ in the limit $n \to \infty$, in the sense that deviations from $\frac{1}{2}$ (e.g., a long run of heads or of tails) will become increasingly rare.

100,000 tosses of a fair (?) coin



José Figueroa-O'Farrill mi4a (Probability) Lecture 9 18 / 23

3 X X 3 X 3

-

Theorem (The (weak) law of large numbers)

Let $X_1, X_2, ...$ be i.i.d. discrete random variables with mean μ and variance σ^2 and let $Z_n = \frac{1}{n}(X_1 + \cdots + X_n)$.

Theorem (The (weak) law of large numbers)

Let $X_1, X_2, ...$ be i.i.d. discrete random variables with mean μ and variance σ^2 and let $Z_n = \frac{1}{n}(X_1 + \cdots + X_n)$. Then

 $\forall \epsilon > 0 \qquad \mathbb{P}(|\mathsf{Z}_n - \mu| < \epsilon) \to 1 \quad \textit{as } n \to \infty$

Theorem (The (weak) law of large numbers)

Let $X_1, X_2, ...$ be i.i.d. discrete random variables with mean μ and variance σ^2 and let $Z_n = \frac{1}{n}(X_1 + \cdots + X_n)$. Then

 $\forall \epsilon > 0 \qquad \mathbb{P}(|Z_n - \mu| < \epsilon) \to 1 \quad \textit{as } n \to \infty$

Proof.

By linearity of expectation, $\mathbb{E}(Z_n) = \mu$,

Theorem (The (weak) law of large numbers)

Let $X_1, X_2, ...$ be i.i.d. discrete random variables with mean μ and variance σ^2 and let $Z_n = \frac{1}{n}(X_1 + \cdots + X_n)$. Then

 $\forall \epsilon > 0 \qquad \mathbb{P}(|Z_n - \mu| < \epsilon) \to 1 \quad \textit{as } n \to \infty$

Proof.

By linearity of expectation, $\mathbb{E}(Z_n) = \mu$, and since the X_i are independent $Var(Z_n) = \frac{1}{n^2} Var(X_1 + \dots + X_n) = \frac{\sigma^2}{n}$.

Theorem (The (weak) law of large numbers)

Let $X_1, X_2, ...$ be i.i.d. discrete random variables with mean μ and variance σ^2 and let $Z_n = \frac{1}{n}(X_1 + \cdots + X_n)$. Then

 $\forall \epsilon > 0 \qquad \mathbb{P}(|Z_n - \mu| < \epsilon) \to 1 \quad \textit{as } n \to \infty$

Proof.

By linearity of expectation, $\mathbb{E}(Z_n) = \mu$, and since the X_i are independent $Var(Z_n) = \frac{1}{n^2} Var(X_1 + \dots + X_n) = \frac{\sigma^2}{n}$. By Chebyshev,

$$\mathbb{P}(|\mathsf{Z}_n - \mu| \ge \varepsilon) \leqslant \frac{\sigma^2}{n\varepsilon^2} \implies \mathbb{P}(|\mathsf{Z}_n - \mu| < \varepsilon) \ge 1 - \frac{\sigma^2}{n\varepsilon^2}$$

We will now justify probability as relative frequency.

20 / 23

A ►

We will now justify probability as relative frequency.

Let $(\Omega, \mathfrak{F}, \mathbb{P})$ be a probability space and let $A \in \mathfrak{F}$ be an event.

We will now justify probability as relative frequency.

Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space and let $A \in \mathcal{F}$ be an event. Let I_A denote the **indicator variable** of A, a discrete random variable defined by

$$I_{A}(\omega) = \begin{cases} 1, & \omega \in A \\ 0, & \omega \notin A \end{cases}$$

We will now justify probability as relative frequency.

Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space and let $A \in \mathcal{F}$ be an event. Let I_A denote the **indicator variable** of A, a discrete random variable defined by

$$I_{A}(\omega) = \begin{cases} 1, & \omega \in A \\ 0, & \omega \notin A \end{cases}$$

The probability mass function f of an indicator variable is very simple: $f(1) = \mathbb{P}(A)$ and hence $f(0) = 1 - \mathbb{P}(A)$.

We will now justify probability as relative frequency.

Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space and let $A \in \mathcal{F}$ be an event. Let I_A denote the **indicator variable** of A, a discrete random variable defined by

$$I_{A}(\omega) = \begin{cases} 1, & \omega \in A \\ 0, & \omega \notin A \end{cases}$$

The probability mass function f of an indicator variable is very simple: $f(1) = \mathbb{P}(A)$ and hence $f(0) = 1 - \mathbb{P}(A)$. Its mean is given by

$$\boldsymbol{\mu} = \mathbb{E}(\mathrm{I}_A) = \boldsymbol{0} \times f(\boldsymbol{0}) + \boldsymbol{1} \times f(\boldsymbol{1}) = f(\boldsymbol{1}) = \mathbb{P}(A)$$

▲□ ▶ ▲ ■ ▶ ▲ ■ ▶ ■ ■ ■ ● ● ●

We will now justify probability as relative frequency.

Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space and let $A \in \mathcal{F}$ be an event. Let I_A denote the **indicator variable** of A, a discrete random variable defined by

$$I_{A}(\omega) = \begin{cases} 1, & \omega \in A \\ 0, & \omega \notin A \end{cases}$$

The probability mass function f of an indicator variable is very simple: $f(1) = \mathbb{P}(A)$ and hence $f(0) = 1 - \mathbb{P}(A)$. Its mean is given by

$$\mu = \mathbb{E}(I_A) = \mathbf{0} \times f(\mathbf{0}) + \mathbf{1} \times f(\mathbf{1}) = f(\mathbf{1}) = \mathbb{P}(A)$$

and its variance by

$$\sigma^{2} = \operatorname{Var}(I_{A}) = (0 - \mu)^{2} f(0) + (1 - \mu)^{2} f(1)$$
$$= \mathbb{P}(A)^{2} (1 - \mathbb{P}(A)) + (1 - \mathbb{P}(A))^{2} \mathbb{P}(A)$$
$$= \mathbb{P}(A)(1 - \mathbb{P}(A))$$

• Now imagine repeating the experiment and counting how many outcomes belong to A.

- Now imagine repeating the experiment and counting how many outcomes belong to A.
- Let X_i denote the random variable which agrees with the indicator variable of A at the ith trial.

- Now imagine repeating the experiment and counting how many outcomes belong to A.
- Let X_i denote the random variable which agrees with the indicator variable of A at the ith trial.
- Then the X₁, X₂,... are i.i.d. discrete random variables, with mean P(A) and variance P(A)(1 − P(A)).

- Now imagine repeating the experiment and counting how many outcomes belong to A.
- Let X_i denote the random variable which agrees with the indicator variable of A at the ith trial.
- Then the X₁, X₂,... are i.i.d. discrete random variables, with mean P(A) and variance P(A)(1 − P(A)).
- Let $Z_n = \frac{1}{n}(X_1 + \dots + X_n)$. What does Z_n measure?

- Now imagine repeating the experiment and counting how many outcomes belong to A.
- Let X_i denote the random variable which agrees with the indicator variable of A at the ith trial.
- Then the X₁, X₂,... are i.i.d. discrete random variables, with mean P(A) and variance P(A)(1 − P(A)).
- Let $Z_n = \frac{1}{n}(X_1 + \cdots + X_n)$. What does Z_n measure?
- Z_n measures the proportion of trials with outcomes in A after n trials. This is what we had originally called N(A)/n.

- Now imagine repeating the experiment and counting how many outcomes belong to A.
- Let X_i denote the random variable which agrees with the indicator variable of A at the ith trial.
- Then the X₁, X₂, ... are i.i.d. discrete random variables, with mean P(A) and variance P(A)(1 − P(A)).
- Let $Z_n = \frac{1}{n}(X_1 + \cdots + X_n)$. What does Z_n measure?
- Z_n measures the proportion of trials with outcomes in A after n trials. This is what we had originally called N(A)/n.
- The law of large numbers says that in the limit as $n \to \infty$, $Z_n \to \mathbb{P}(A)$ in probability.

- Now imagine repeating the experiment and counting how many outcomes belong to A.
- Let X_i denote the random variable which agrees with the indicator variable of A at the ith trial.
- Then the X₁, X₂, ... are i.i.d. discrete random variables, with mean P(A) and variance P(A)(1 − P(A)).
- Let $Z_n = \frac{1}{n}(X_1 + \cdots + X_n)$. What does Z_n measure?
- Z_n measures the proportion of trials with outcomes in A after n trials. This is what we had originally called N(A)/n.
- The law of large numbers says that in the limit as $n \to \infty$, $Z_n \to \mathbb{P}(A)$ in probability.
- This makes precise our initial hand-waving argument of N(A)/n "converging in some way" to P(A).

• X, Y independent: $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$

□ > < E > < E > E = のへで

- X, Y independent: $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$
- $\mathbb{E}(XY)$ defines an inner product

- X, Y independent: $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$
- $\mathbb{E}(XY)$ defines an inner product
- X, Y independent: Var(X + Y) = Var(X) + Var(Y)

- X, Y independent: $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$
- $\mathbb{E}(XY)$ defines an inner product
- X, Y independent: Var(X + Y) = Var(X) + Var(Y)
- In general: Var(X + Y) = Var(X) + Var(Y) + 2 Cov(X, Y)

- X, Y independent: $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$
- $\mathbb{E}(XY)$ defines an inner product
- X, Y independent: Var(X + Y) = Var(X) + Var(Y)
- In general: Var(X + Y) = Var(X) + Var(Y) + 2 Cov(X, Y)
- covariance: $Cov(X, Y) = \mathbb{E}(XY) \mathbb{E}(X)\mathbb{E}(Y)$. If

Cov(X, Y) = 0 we say X, Y are uncorrelated

▲□ ▶ ▲ □ ▶ ▲ □ ▶ □ □ ● ● ● ●

- X, Y independent: $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$
- $\mathbb{E}(XY)$ defines an inner product
- X, Y independent: Var(X + Y) = Var(X) + Var(Y)
- In general: Var(X + Y) = Var(X) + Var(Y) + 2 Cov(X, Y)
- covariance: Cov(X, Y) = E(XY) − E(X)E(Y). If Cov(X, Y) = 0 we say X, Y are uncorrelated
- correlation: $\rho(X, Y) = Cov(X, Y)/(\sigma(X)\sigma(Y))$ measures "linear dependence" between X, Y

- X, Y independent: $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$
- $\mathbb{E}(XY)$ defines an inner product
- X, Y independent: Var(X + Y) = Var(X) + Var(Y)
- In general: Var(X + Y) = Var(X) + Var(Y) + 2 Cov(X, Y)
- covariance: Cov(X, Y) = E(XY) E(X)E(Y). If Cov(X, Y) = 0 we say X, Y are uncorrelated
- correlation: $\rho(X, Y) = Cov(X, Y)/(\sigma(X)\sigma(Y))$ measures "linear dependence" between X, Y
- We proved two inequalities:

- X, Y independent: $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$
- $\mathbb{E}(XY)$ defines an inner product
- X, Y independent: Var(X + Y) = Var(X) + Var(Y)
- In general: Var(X + Y) = Var(X) + Var(Y) + 2 Cov(X, Y)
- covariance: Cov(X, Y) = E(XY) E(X)E(Y). If Cov(X, Y) = 0 we say X, Y are uncorrelated
- correlation: $\rho(X, Y) = Cov(X, Y)/(\sigma(X)\sigma(Y))$ measures "linear dependence" between X, Y
- We proved two inequalities:
 - Markov: $\mathbb{P}(|X| \ge \alpha) \le \mathbb{E}(|X|)/\alpha$

- X, Y independent: $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$
- $\mathbb{E}(XY)$ defines an inner product
- X, Y independent: Var(X + Y) = Var(X) + Var(Y)
- In general: Var(X + Y) = Var(X) + Var(Y) + 2 Cov(X, Y)
- covariance: Cov(X, Y) = E(XY) E(X)E(Y). If Cov(X, Y) = 0 we say X, Y are uncorrelated
- correlation: $\rho(X, Y) = Cov(X, Y)/(\sigma(X)\sigma(Y))$ measures "linear dependence" between X, Y
- We proved two inequalities:
 - Markov: $\mathbb{P}(|X| \ge \alpha) \le \mathbb{E}(|X|)/\alpha$
 - Chebyshev: $\mathbb{P}(|X-\mu| \geqslant \epsilon) \leqslant \sigma^2/\epsilon^2$

- X, Y independent: $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$
- $\mathbb{E}(XY)$ defines an inner product
- X, Y independent: Var(X + Y) = Var(X) + Var(Y)
- In general: Var(X + Y) = Var(X) + Var(Y) + 2 Cov(X, Y)
- covariance: Cov(X, Y) = E(XY) E(X)E(Y). If Cov(X, Y) = 0 we say X, Y are uncorrelated
- correlation: $\rho(X, Y) = Cov(X, Y)/(\sigma(X)\sigma(Y))$ measures "linear dependence" between X, Y
- We proved two inequalities:
 - Markov: $\mathbb{P}(|X| \ge \alpha) \le \mathbb{E}(|X|)/\alpha$
 - Chebyshev: $\mathbb{P}(|X-\mu| \geqslant \epsilon) \leqslant \sigma^2/\epsilon^2$
- The **law of large numbers** "explains" the relative frequency definition of probability:

- X, Y independent: $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$
- $\mathbb{E}(XY)$ defines an inner product
- X, Y independent: Var(X + Y) = Var(X) + Var(Y)
- In general: Var(X + Y) = Var(X) + Var(Y) + 2 Cov(X, Y)
- covariance: Cov(X, Y) = E(XY) E(X)E(Y). If Cov(X, Y) = 0 we say X, Y are uncorrelated
- correlation: $\rho(X, Y) = Cov(X, Y)/(\sigma(X)\sigma(Y))$ measures "linear dependence" between X, Y
- We proved two inequalities:
 - Markov: $\mathbb{P}(|X| \ge \alpha) \le \mathbb{E}(|X|)/\alpha$
 - Chebyshev: $\mathbb{P}(|X-\mu| \geqslant \epsilon) \leqslant \sigma^2/\epsilon^2$
- The law of large numbers "explains" the relative frequency definition of probability: it says that if X_i are i.i.d. discrete random variables, then as n → ∞, ¹/_n(X₁ + ··· + X_n) → μ in probability;

- X, Y independent: $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$
- $\mathbb{E}(XY)$ defines an inner product
- X, Y independent: Var(X + Y) = Var(X) + Var(Y)
- In general: Var(X + Y) = Var(X) + Var(Y) + 2 Cov(X, Y)
- covariance: Cov(X, Y) = E(XY) E(X)E(Y). If Cov(X, Y) = 0 we say X, Y are uncorrelated
- correlation: $\rho(X, Y) = Cov(X, Y)/(\sigma(X)\sigma(Y))$ measures "linear dependence" between X, Y
- We proved two inequalities:
 - Markov: $\mathbb{P}(|X| \ge \alpha) \le \mathbb{E}(|X|)/\alpha$
 - Chebyshev: $\mathbb{P}(|X-\mu| \geqslant \epsilon) \leqslant \sigma^2/\epsilon^2$
- The **law of large numbers** "explains" the relative frequency definition of probability: it says that if X_i are i.i.d. discrete random variables, then as $n \to \infty$, $\frac{1}{n}(X_1 + \dots + X_n) \to \mu$ *in probability*; i.e., deviations from μ are still possible, but they are increasingly improbable

Proof of the Cauchy–Schwarz inequality

The Cauchy–Schwarz inequality says that if x, y are any two vectors in a positive-definite inner product space, then

 $|\langle x,y\rangle|\leqslant |x||y|\;,\qquad \text{where }|x|=\sqrt{\langle x,x\rangle}\;\text{is the length}.$

▲■▶ ▲ ■▶ ★ ■▶ 三日= のへで

Proof of the Cauchy–Schwarz inequality

The Cauchy–Schwarz inequality says that if x, y are any two vectors in a positive-definite inner product space, then

 $|\left\langle x,y\right\rangle |\leqslant |x||y|\;,\qquad \text{where }|x|=\sqrt{\left\langle x,x\right\rangle }\text{ is the length.}$

Any two vectors lie on a plane, so let's pretend we are in \mathbb{R}^2 ,

◆◎ ▶ ◆ ■ ▶ ◆ ■ ▶ ● ● ● ● ● ● ●

Proof of the Cauchy–Schwarz inequality

The Cauchy–Schwarz inequality says that if x, y are any two vectors in a positive-definite inner product space, then

 $|\langle x,y \rangle| \leqslant |x||y|$, where $|x| = \sqrt{\langle x,x \rangle}$ is the length.

Any two vectors lie on a plane, so let's pretend we are in \mathbb{R}^2 , and diagonalising $\langle -, - \rangle$, we take it to be the dot product.

◆◎ ▶ ◆ ■ ▶ ◆ ■ ▶ ● ● ● ● ● ● ●
Proof of the Cauchy–Schwarz inequality

The Cauchy–Schwarz inequality says that if x, y are any two vectors in a positive-definite inner product space, then

 $|\langle x,y \rangle| \leqslant |x||y|$, where $|x| = \sqrt{\langle x,x \rangle}$ is the length.

Any two vectors lie on a plane, so let's pretend we are in \mathbb{R}^2 , and diagonalising $\langle -, - \rangle$, we take it to be the dot product. In that case,

 $\mathbf{x} \cdot \mathbf{y} = |\mathbf{x}| |\mathbf{y}| \cos \theta$,

where θ is the angle between x and y. Since $|\cos \theta| \le 1$, the inequality follows.



Back to the main story.