Introduction to Statistics Session 2

Grzegorz Chrupała

Saarland University

October 13, 2011

3

-

• • • • • • • • • • •

Random variables

- Function $X: \Omega \to \mathbb{R}^n$ (typically n = 1)
- It may be more convenient to work with real number than directly with events
- Coin toss: $X : \{H, T\} \rightarrow \{0, 1\}$
- Sum of two dice throws: $\{1..6\}^2 \rightarrow \{2..12\}$
- Probability mass function:

$$\mathrm{p}(x) = P(X = x) = P(A)$$
 where $A = \{\omega \in \Omega : X(\omega) = x\}$

・ 同 ト ・ ヨ ト ・ ヨ ト

Expectation

• Expectation is a mean (weighted average) of a random variable

$$E(X) = \sum_{x} p(x) \cdot x$$

• Example: rolling a dice:

$$E(X) =$$

3

(日) (同) (三) (三)

Expectation

• Expectation is a mean (weighted average) of a random variable

$$E(X) = \sum_{x} p(x) \cdot x$$

• Example: rolling a dice:

$$E(X) = \sum_{x=1}^{6} p(x)x = \sum_{x=1}^{6} \frac{x}{6} = 3.5$$

• A function g(X) defines new random variable. In this case:

$$E(g(X)) = \sum_{x} p(x)g(x)$$

Example?

• Also for two random variables:

$$E(X+Y)=E(X)+E(Y)$$

and if independent

$$E(XY) = E(X)E(Y) \rightarrow A = A$$

Chrupala (Saarland)

Variance

• Variance measures how much values of a random variable vary

$$Var(X) = E[(X - E(X))^2]$$

- $\bullet\,$ Standard deviation σ is the square root of the variance
- What is the variance of a random variable describing a single throw of a dice?



- Entropy is a measure of degree of uncertainty.
- The most important concept in information theory
- Entropy is a property of a random variable X distributed according the pmf p

$$H(X) = H(p) = E(-\log_2(x)) = -\sum_x p(x) \log_2(p(x))$$

• For $\log_2(x)$ units are bits, for $\ln(x)$, nats

Entropy as amount of information

- You can think of entropy as measuring the cost of transmitting information about the result of an experiment
- Fair coin toss:

Entropy as amount of information

- You can think of entropy as measuring the cost of transmitting information about the result of an experiment
- Fair coin toss:

$$H(X) = -\sum_{x=0}^{1} p(x) \log_2(p(x))$$
(1)
= $\frac{1}{2} [-\log_2\left(\frac{1}{2}\right) - \log_2\left(\frac{1}{2}\right)]$ (2)
= $\frac{1}{2} \cdot 2$ (3)

Entropy of an unfair coin



Properties of entropy

- *H*(*p*) ≥ 0
- When is it H(p) = 0?
- The highest entropy corresponds to the most uniform distribution

3

Entropy: joint and conditional

• For two variables X and Y, the amount of information needed to specify values of both

$$H(X,Y) = -\sum_{x}\sum_{y} p(x,y) \log_2(p(x,y))$$

• Conditional entropy: if we know the value of X, how much does to cost to transmit the value of Y?

$$H(Y|X) = \sum_{x} p(x)H(Y|X=x)$$
(4)

$$= \sum_{x} p(x) \left[-\sum_{y} p(y|x) \log(p(y|x)) \right]$$
(5)

$$= -\sum_{x}\sum_{y} p(y|x)p(x)\log(p(y|x))$$
(6)

$$= -\sum_{x,y} p(x,y) \log(p(y|x))$$
(7)

Conditional entropy: example

- X: number of heads in two tosses of a fair coin
- Y: is at least one of two tosses of a fair coin a heads?

	Х	Υ
HH	2	1
ΗT	1	1
ΤT	0	0
ΤH	1	1

What is H(X)? What is H(X|Y)?

Chain rule for entropy

$$H(X, Y) = H(X|Y) + H(Y)$$

$$H(X_1, ..., X_n) = H(X_1) + H(X_2|X_1) + \dots + H(X_n|X_1, ..., X_{n-1})$$

- ∢ ≣ →

Image: A image: A

2

Mutual information

• From the chain rule we have

$$H(X) + H(Y|X) = H(Y) + H(X|Y)$$

Therefore

$$H(X) - H(X|Y) = H(Y) - H(Y|X)$$

• This difference is known as Mutual information I(X; Y)

Image: A matrix of the second seco

Joint and conditional entropy and mutual information



3

(人間) トイヨト イヨト

Mutual information

$$I(X; Y) = H(X) - H(X|Y)$$

= $H(X) + H(Y) + H(X, Y)$
...
= $\sum_{x} \sum_{y} p(x, y) \log \left(\frac{p(x, y)}{p(x)p(y)}\right)$

.

• What is I(X; X)?

æ

イロト イ団ト イヨト イヨト

Kullback Leibler divergence

• A measure of the difference between two probability mass functions p and q is Kullback Leibler divergence (relative entropy)

$$D(p||q) = \sum_{x} \mathrm{p}(x) \log\left(rac{\mathrm{p}(x)}{\mathrm{q}(x)}
ight)$$

- $\bullet\,$ Can be interpreted as an average number of bits wasted by encoding events distributed according to p with a code based on q
- We can define mutual information in terms of KL divergence:

$$I(X; Y) = D(p(x, y)||p(x)p(y))$$