## COMP2610/COMP6261 - Information Theory Lecture 9: Probabilistic Inequalities

#### Mark Reid and Aditya Menon

Research School of Computer Science The Australian National University



### August 19th, 2014

Mutual information chain rule

Jensen's inequality

"Information cannot hurt"

Data processing inequality

### Theorem

if 
$$X \to Y \to Z$$
 then:  $I(X; Y) \ge I(X; Z)$ 

- X is the state of the world, Y is the data gathered and Z is the processed data
- No "clever" manipulation of the data can improve the inferences that can be made from the data
- No processing of Y, deterministic or random, can increase the information that Y contains about X

- Markov's inequality
- Chebyshev's inequality
- Law of large numbers

- 1 Properties of expectation and variance
- 2 Markov's inequality
- 3 Chebyshev's inequality
- 4 Law of large numbers

## 5 Wrapping Up

### 1 Properties of expectation and variance

### 2 Markov's inequality

- 3 Chebyshev's inequality
- 4 Law of large numbers

## 5 Wrapping Up

Let X be a random variable over  $\mathcal{X}$ , with probability distribution p

Expected value:

$$\mathbb{E}[X] = \sum_{x \in \mathcal{X}} x \cdot p(x).$$

Variance:

$$\mathbb{V}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2]$$
  
=  $\mathbb{E}[X^2] - (\mathbb{E}[X])^2.$ 

Standard deviation is  $\sqrt{\mathbb{V}[X]}$ 

A key property of expectations is linearity:

$$\mathbb{E}\left[\sum_{i=1}^{n} X_{i}\right] = \sum_{i=1}^{n} \mathbb{E}\left[X_{i}\right].$$

This holds even if the variables are dependent!

We have for any  $a \in \mathbb{R}$ ,

$$\mathbb{E}[aX] = a \cdot \mathbb{E}[X].$$

We have linearity of variance for independent random variables:

$$\mathbb{V}\left[\sum_{i=1}^{n} X_{i}\right] = \sum_{i=1}^{n} \mathbb{V}\left[X_{i}\right].$$

Does not hold if the variables are dependent

We have for any  $a \in \mathbb{R}$ ,

$$\mathbb{V}[aX] = a^2 \cdot \mathbb{V}[X].$$

Properties of expectation and variance

2 Markov's inequality

3 Chebyshev's inequality

4 Law of large numbers

5 Wrapping Up

1000 school students sit an examination

The busy principal is only told that the average score is 40 out of 100

The principal wants to estimate the maximum number of students who scored more than 80

• Would it make sense to ask about the minimum number of students?

Call x the number of students who score more than 80

We know:

$$40 \cdot 1000 = 80x + S$$
,

where S is the total score of students scoring less than 80

Exam scores are nonnegative, so certainly  $S \ge 0$ 

Thus,  $80x \le 40 \cdot 1000$ , or

 $x \leq 500.$ 

Can we formalise this more generally?

#### Theorem

Let X be a nonnegative random variable. Then, for any  $\lambda > 0$ ,

$$p(X \ge \lambda) \le \frac{\mathbb{E}[X]}{\lambda}.$$

Bounds probability of observing a large outcome

#### Corollary

Let X be a nonnegative random variable. Then, for any  $\lambda > 0$ ,

$$p(X \ge \lambda \cdot \mathbb{E}[X]) \le \frac{1}{\lambda}.$$

Observations of nonnegative random variable unlikely to be much larger than expected value

$$\mathbb{E}[X] = \sum_{x \in \mathcal{X}} x \cdot p(x)$$
  
 $= \sum_{x < \lambda} x \cdot p(x) + \sum_{x \ge \lambda} x \cdot p(x)$   
 $\ge \sum_{x \ge \lambda} x \cdot p(x)$  nonneg. of random variable  
 $\ge \sum_{x \ge \lambda} \lambda \cdot p(x)$   
 $= \lambda \cdot p(X \ge \lambda).$ 

## Markov's Inequality

Illustration from http://justindomke.wordpress.com/



## Markov's Inequality

Illustration from http://justindomke.wordpress.com/



## Markov's Inequality

Illustration from http://justindomke.wordpress.com/



Mark Reid and Aditya Menon (ANU) COMP2610/COMP6261 - Information Theory

## Markov's Inequality Illustration from http://justindomke.wordpress.com/



Properties of expectation and variance

2 Markov's inequality

3 Chebyshev's inequality

4 Law of large numbers

### 5 Wrapping Up

# Chebyshev's Inequality Motivation

Markov's inequality only uses the mean of the distribution

What about the spread of the distribution (variance)?



#### Theorem

Let X be a random variable with  $\mathbb{E}[X] < \infty$ . Then, for any  $\lambda > 0$ ,

$$p(|X - \mathbb{E}[X]| \ge \lambda) \le \frac{\mathbb{V}[X]}{\lambda^2}.$$

Bounds the probability of observing an "unexpected" outcome

Do not require non negativity

Two-sided bound

#### Corollary

Let X be a random variable with  $\mathbb{E}[X] < \infty$ . Then, for any  $\lambda > 0$ ,

$$p(|X - \mathbb{E}[X]| \ge \lambda \cdot \sqrt{\mathbb{V}[X]}) \le rac{1}{\lambda^2}.$$

Observations are unlikely to occur several standard deviations away from the mean

# Chebyshev's Inequality Proof

Define

$$Y = (X - \mathbb{E}[X])^2.$$

Then, by Markov's inequality, for any  $\nu > 0$ ,

$$p(Y \ge \nu) \le \frac{\mathbb{E}[Y]}{\nu}.$$

But,

$$\mathbb{E}[Y] = \mathbb{V}[X].$$

Also,

$$Y \ge \nu \iff |X - \mathbb{E}[X]| \ge \sqrt{\nu}.$$

Thus, setting  $\lambda = \sqrt{\nu}$ ,

$$p(|X - \mathbb{E}[X]| \ge \lambda) \le \frac{\mathbb{V}[X]}{\lambda^2}.$$

## Chebyshev's Inequality

Illustration

For a binomial with N trials and success probability  $\theta$ , we have e.g.

$$p(|X - N\theta| \ge \sqrt{2N\theta(1-\theta)}) \le \frac{1}{2}$$



Suppose we have a coin with bias  $\theta$ , i.e.  $p(X = 1) = \theta$ 

Say we flip the coin *n* times, and observe  $x_1, \ldots, x_n \in \{0, 1\}$ 

We use the maximum likelihood estimator of  $\theta$ :

$$\hat{\theta}_n = \frac{x_1 + \ldots + x_n}{n}$$

Estimate how large *n* should be such that

$$p(|\hat{ heta}_n - heta| \ge 0.05) \le 0.01?$$

1% probability of a 5% error

# Chebyshev's Inequality Example

Observe that

$$\mathbb{E}[\hat{\theta}_n] = \frac{\sum_{i=1}^n \mathbb{E}[x_i]}{n} = \theta$$
$$\mathbb{V}[\hat{\theta}_n] = \frac{\sum_{i=1}^n \mathbb{V}[x_i]}{n^2} = \frac{\theta(1-\theta)}{n}.$$

Thus, applying Chebyshev's inequality to  $\hat{\theta}_n$ ,

$$p(|\hat{ heta}_n - heta| > 0.05) \leq rac{ heta(1 - heta)}{(0.05)^2 \cdot n}$$

We are guaranteed this is less than 0.01 if

$$n \geq \frac{\theta(1- heta)}{(0.05)^2(0.01)}.$$

When  $\theta = 0.5$ ,  $n \ge 10,000!$ 

Mark Reid and Aditya Menon (ANU) COMP2610/COMP6261 - Information Theory

Properties of expectation and variance

2 Markov's inequality

3 Chebyshev's inequality

4 Law of large numbers

### 5 Wrapping Up

Let  $X_1, \ldots, X_n$  be random variables such that:

- Each X<sub>i</sub> is independent of X<sub>j</sub>
- The distribution of X<sub>i</sub> is the same as that of X<sub>j</sub>

Then, we say that  $X_1, \ldots, X_n$  are independent and identically distributed (or iid)

Example: For *n* independent flips of an unbiased coin,  $X_1, \ldots, X_n$  are iid from Bernoulli $(\frac{1}{2})$ 

## Law of Large Numbers

#### Theorem

Let  $X_1, \ldots, X_n$  be a sequence of iid random variables, with

$$\mathbb{E}[X_i] = \mu$$

and  $\mathbb{V}[X_i] < \infty$ . Define

$$\bar{X}_n = \frac{X_1 + \ldots + X_n}{n}$$

Then, for any  $\epsilon > 0$ ,

$$\lim_{n\to\infty}p(|\bar{X}_n-\mu|>\epsilon)=0.$$

Given enough trials, the empirical "success frequency" will be close to the expected value

# Law of Large Numbers Proof

Since  $X_i$ 's are identically distributed,

$$\mathbb{E}[\bar{X}_n] = \mu.$$

Since the  $X_i$ 's are independent,

$$\mathbb{V}[\bar{X}_n] = \mathbb{V}\left[\frac{X_1 + \ldots + X_n}{n}\right]$$
$$= \frac{\mathbb{V}\left[X_1 + \ldots + X_n\right]}{n^2}$$
$$= \frac{n\sigma^2}{n^2}$$
$$= \frac{\sigma^2}{n}.$$

Applying Chebyshev's inequality to  $\bar{X}_n$ ,

$$p(|\bar{X}_n - \mu| \ge \epsilon) \le \frac{\mathbb{V}[\bar{X}_n]}{\epsilon^2} = \frac{\sigma^2}{n\epsilon^2}.$$

As  $n \to \infty$ , the right hand side  $\to 0$ .

Thus,

$$p(|\bar{X}_n - \mu| < \epsilon) \to 1.$$

# Law of Large Numbers

N = 1000 trials with Bernoulli random variable with parameter  $\frac{1}{2}$ 



## Law of Large Numbers

N = 50000 trials with Bernoulli random variable with parameter  $\frac{1}{2}$ 



Properties of expectation and variance

2 Markov's inequality

3 Chebyshev's inequality

4 Law of large numbers



- Markov's inequality
- Chebyshev's inequality
- Law of large numbers

- Ensembles and sequences
- Typical sets
- Approximation Equipartition (AEP)