

Note to other teachers and users of these slides. Andrew would be delighted if you found this source material useful in giving your own lectures. Feel free to use these slides verbatim, or to modify them to fit your own needs. PowerPoint originals are available. If you make use of a significant portion of these slides in your own lecture, please include this message, or the following link to the source repository of Andrew's tutorials: <http://www.cs.cmu.edu/~awm/tutorials> . Comments and corrections gratefully received.

# Information Gain

**Andrew W. Moore**  
**Professor**  
**School of Computer Science**  
**Carnegie Mellon University**

[www.cs.cmu.edu/~awm](http://www.cs.cmu.edu/~awm)

[awm@cs.cmu.edu](mailto:awm@cs.cmu.edu)

412-268-7599

# Bits

You are watching a set of independent random samples of  $X$

You see that  $X$  has four possible values

$P(X=A) = 1/4$	$P(X=B) = 1/4$	$P(X=C) = 1/4$	$P(X=D) = 1/4$
----------------	----------------	----------------	----------------

So you might see: BAACBADCDADDDA...

You transmit data over a binary serial link. You can encode each reading with two bits (e.g.  $A = 00$ ,  $B = 01$ ,  $C = 10$ ,  $D = 11$ )

0100001001001110110011111100...

# Fewer Bits

Someone tells you that the probabilities are not equal

$P(X=A) = 1/2$	$P(X=B) = 1/4$	$P(X=C) = 1/8$	$P(X=D) = 1/8$
----------------	----------------	----------------	----------------

It's possible...

...to invent a coding for your transmission that only uses 1.75 bits on average per symbol. How?

# Fewer Bits

Someone tells you that the probabilities are not equal

$P(X=A) = 1/2$	$P(X=B) = 1/4$	$P(X=C) = 1/8$	$P(X=D) = 1/8$
----------------	----------------	----------------	----------------

It's possible...

...to invent a coding for your transmission that only uses 1.75 bits on average per symbol. How?

A	0
B	10
C	110
D	111

(This is just one of several ways)

# Fewer Bits

Suppose there are three equally likely values...

$P(X=A) = 1/3$	$P(X=B) = 1/3$	$P(X=C) = 1/3$
----------------	----------------	----------------

Here's a naïve coding, costing 2 bits per symbol

A	00
B	01
C	10

Can you think of a coding that would need only 1.6 bits per symbol on average?

In theory, it can in fact be done with 1.58496 bits per symbol.

# General Case

Suppose  $X$  can have one of  $m$  values...  $V_1, V_2, \dots, V_m$

$P(X=V_1) = p_1$	$P(X=V_2) = p_2$	....	$P(X=V_m) = p_m$
------------------	------------------	------	------------------

What's the smallest possible number of bits, on average, per symbol, needed to transmit a stream of symbols drawn from  $X$ 's distribution? It's

$$\begin{aligned} H(X) &= -p_1 \log_2 p_1 - p_2 \log_2 p_2 - \dots - p_m \log_2 p_m \\ &= -\sum_{j=1}^m p_j \log_2 p_j \end{aligned}$$

$H(X)$  = The entropy of  $X$

- "High Entropy" means  $X$  is from a uniform (boring) distribution
- "Low Entropy" means  $X$  is from varied (peaks and valleys) distribution

# General Case

Suppose X can have one of  $m$  values...  $V_1, V_2, \dots, V_m$

$P(X=V_1) = p_1$	$P(X=V_2) = p_2$	....	$P(X=V_m) = p_m$
------------------	------------------	------	------------------

What's the smallest possible number of symbols in a stream of length  $n$  that contains  $nH(X)$  symbols of X's distribution?

A histogram of the frequency distribution of values of X would have many lows and one or two highs

A histogram of the frequency distribution of values of X would be flat

$$H(X) = -\sum_{j=1}^m p_j \log_2 p_j$$

$H(X)$  = The entropy of X

- "High Entropy" means X is from a uniform (boring) distribution
- "Low Entropy" means X is from varied (peaks and valleys) distribution

# General Case

Suppose  $X$  can have one of  $m$  values...  $V_1, V_2, \dots, V_m$

$P(X=V_1) = p_1$	$P(X=V_2) = p_2$	....	$P(X=V_m) = p_m$
------------------	------------------	------	------------------

What's the smallest possible number of symbols in a stream of values of  $X$  that would contain all the information in  $X$ 's distribution?

A histogram of the frequency distribution of values of  $X$  would be flat

A histogram of the frequency distribution of values of  $X$  would have many lows and one or two highs

..and so the values sampled from it would be all over the place

..and so the values sampled from it would be more predictable

$$H(X) = -\sum_{i=1}^m p_i \log_2 p_i$$

$$= -\sum_{i=1}^m p_i \log_2 p_i$$

$H(X)$  = The entropy of  $X$

- "High Entropy" means  $X$  is from a uniform (boring) distribution
- "Low Entropy" means  $X$  is from varied (peaks and valleys) distribution

# Entropy in a nut-shell



Low Entropy



High Entropy

# Entropy in a nut-shell



Low Entropy

..the values (locations of soup) sampled entirely from within the soup bowl



High Entropy

..the values (locations of soup) unpredictable... almost uniformly sampled throughout our dining room

# Specific Conditional Entropy $H(Y|X=v)$

**Suppose I'm trying to predict output Y and I have input X**

**X = College Major**

**Y = Likes "Gladiator"**

X	Y
Math	Yes
History	No
CS	Yes
Math	No
Math	No
CS	Yes
History	No
Math	Yes

**Let's assume this reflects the true probabilities**

**E.G. From this data we estimate**

- $P(\text{LikeG} = \text{Yes}) = 0.5$
- $P(\text{Major} = \text{Math} \ \& \ \text{LikeG} = \text{No}) = 0.25$
- $P(\text{Major} = \text{Math}) = 0.5$
- $P(\text{LikeG} = \text{Yes} \mid \text{Major} = \text{History}) = 0$

**Note:**

- $H(X) = 1.5$
- $H(Y) = 1$

# Specific Conditional Entropy $H(Y|X=v)$

**X = College Major**

**Y = Likes "Gladiator"**

**Definition of Specific Conditional Entropy:**

$H(Y|X=v)$  = **The entropy of  $Y$  among only those records in which  $X$  has value  $v$**

X	Y
Math	Yes
History	No
CS	Yes
Math	No
Math	No
CS	Yes
History	No
Math	Yes

# Specific Conditional Entropy $H(Y|X=v)$

**X = College Major**

**Y = Likes "Gladiator"**

X	Y
Math	Yes
History	No
CS	Yes
Math	No
Math	No
CS	Yes
History	No
Math	Yes

**Definition of Specific Conditional Entropy:**

$H(Y|X=v)$  = **The entropy of  $Y$  among only those records in which  $X$  has value  $v$**

**Example:**

- $H(Y|X=Math) = 1$
- $H(Y|X=History) = 0$
- $H(Y|X=CS) = 0$

# Conditional Entropy $H(Y|X)$

**X = College Major**

**Y = Likes "Gladiator"**

## Definition of Conditional Entropy:

$H(Y|X)$  = The average specific conditional entropy of  $Y$

= if you choose a record at random what will be the conditional entropy of  $Y$ , conditioned on that row's value of  $X$

= Expected number of bits to transmit  $Y$  if both sides will know the value of  $X$

$$= \sum_j \text{Prob}(X=v_j) H(Y | X = v_j)$$

X	Y
Math	Yes
History	No
CS	Yes
Math	No
Math	No
CS	Yes
History	No
Math	Yes

# Conditional Entropy

**X = College Major**

**Y = Likes "Gladiator"**

**Definition of Conditional Entropy:**

$H(Y|X)$  = The average conditional entropy of  $Y$

$$= \sum_j \text{Prob}(X=v_j) H(Y | X = v_j)$$

**Example:**

$v_j$	$\text{Prob}(X=v_j)$	$H(Y   X = v_j)$
Math	0.5	1
History	0.25	0
CS	0.25	0

$$H(Y|X) = 0.5 * 1 + 0.25 * 0 + 0.25 * 0 = 0.5$$

X	Y
Math	Yes
History	No
CS	Yes
Math	No
Math	No
CS	Yes
History	No
Math	Yes

# Information Gain

**X = College Major**

**Y = Likes "Gladiator"**

**Definition of Information Gain:**

$IG(Y|X)$  = I must transmit  $Y$ .  
How many bits on average  
would it save me if both ends of  
the line knew  $X$ ?

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

**Example:**

- $H(Y) = 1$
- $H(Y|X) = 0.5$
- Thus  $IG(Y|X) = 1 - 0.5 = 0.5$

X	Y
Math	Yes
History	No
CS	Yes
Math	No
Math	No
CS	Yes
History	No
Math	Yes

# Information Gain Example

wealth values: poor rich

gender Female 14423 1769   $H(\text{wealth} | \text{gender} = \text{Female}) = 0.497654$

Male 22732 9918   $H(\text{wealth} | \text{gender} = \text{Male}) = 0.885847$

$H(\text{wealth}) = 0.793844$   $H(\text{wealth} | \text{gender}) = 0.757154$

$IG(\text{wealth} | \text{gender}) = 0.0366896$

# Another example

wealth values: poor rich

agegroup	10s	2507	3		$H(\text{wealth} \mid \text{agegroup} = 10s) = 0.0133271$
	20s	11262	743		$H(\text{wealth} \mid \text{agegroup} = 20s) = 0.334906$
	30s	9468	3461		$H(\text{wealth} \mid \text{agegroup} = 30s) = 0.838134$
	40s	6738	3986		$H(\text{wealth} \mid \text{agegroup} = 40s) = 0.951961$
	50s	4110	2509		$H(\text{wealth} \mid \text{agegroup} = 50s) = 0.957376$
	60s	2245	809		$H(\text{wealth} \mid \text{agegroup} = 60s) = 0.834049$
	70s	668	147		$H(\text{wealth} \mid \text{agegroup} = 70s) = 0.680882$
	80s	115	16		$H(\text{wealth} \mid \text{agegroup} = 80s) = 0.535474$
	90s	42	13		$H(\text{wealth} \mid \text{agegroup} = 90s) = 0.788941$

$H(\text{wealth}) = 0.793844$     $H(\text{wealth} \mid \text{agegroup}) = 0.709463$

$IG(\text{wealth} \mid \text{agegroup}) = 0.0843813$

# Relative Information Gain

**X = College Major**

**Y = Likes "Gladiator"**

**Definition of Relative Information Gain:**

*RIG(Y|X)* = I must transmit *Y*, what fraction of the bits on average would it save me if both ends of the line knew *X*?

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

**Example:**

- $H(Y|X) = 0.5$
- $H(Y) = 1$
- **Thus  $IG(Y|X) = (1 - 0.5)/1 = 0.5$**

X	Y
Math	Yes
History	No
CS	Yes
Math	No
Math	No
CS	Yes
History	No
Math	Yes

# What is Information Gain used for?

Suppose you are trying to predict whether someone is going live past 80 years. From historical data you might find...

- $IG(\text{LongLife} \mid \text{HairColor}) = 0.01$
- $IG(\text{LongLife} \mid \text{Smoker}) = 0.2$
- $IG(\text{LongLife} \mid \text{Gender}) = 0.25$
- $IG(\text{LongLife} \mid \text{LastDigitOfSSN}) = 0.00001$

IG tells you how interesting a 2-d contingency table is going to be.