



# Entropy and decisions

CSC401/2511 – Natural Language Computing – Winter 2022

Lecture 5

University of Toronto

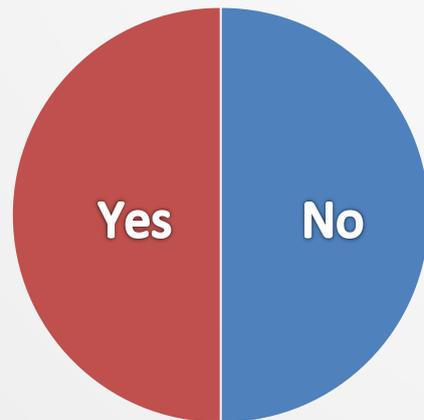
# This lecture

- Information theory.
  - Entropy.
  - Mutual information, etc.
- Decisions.
  - Classification.
  - Significance and hypothesis testing.

*Can we quantify the statistical structure in a model of communication?  
Can we quantify the meaningful difference between statistical models?*

# Information

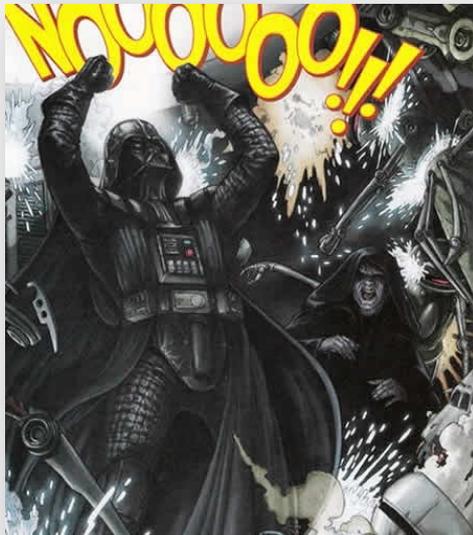
- Imagine Darth Vader is about to say either “yes” or “no” with **equal** probability.
  - You don’t know what he’ll say.
- You have a certain amount of **uncertainty** – a lack of information.



Darth Vader is © Disney  
And the prequels and Rey/Finn Star Wars suck

# Information

- Imagine you then **observe** Darth Vader saying “no”
- Your uncertainty is **gone**; you’ve **received information**.
- **How much** information do you **receive** about event  $x$  when you observe it?



“Choosing 1 out of 2” gives a bit of information

$$I(x) = \mathbf{1 \text{ bit}} \text{ for } P(x) = \frac{1}{2}$$

# Information

- Imagine there is both DARTH Vader and VARTH Dader.
- Observing what both DV and VD say gives us 2 bits of information.
- There are  $2^2$  scenarios with equal possibilities:
  - Yes/Yes, Yes/No, No/Yes, No/No

Darth Vader



Varth Dader



# Information

- So  $I(x)=2$  bits is brought by  $P(x) = \frac{1}{2^2}$
- $I(x)$  doubles when  $\frac{1}{P(x)}$  is squared.
- Let's describe  $I(x)$  with **negative log likelihood**:

$$I(x) = \log_2 \frac{1}{P(x)}$$

For capturing the  
Logarithm relationship

$I(x) = -\log_2 P(x)$ ;  
So here comes the negation

Going back to the “yes/no” example:

$$I(\text{no}) = \log_2 \frac{1}{P(\text{no})} = \log_2 \frac{1}{1/2} = \mathbf{1 \text{ bit}}$$

Note 1: Negative log likelihood is also called **surprisal**.

Note 2: information contents computed with log base 2 has unit “bit”. Log base e => unit “nat”.

# Information

- Imagine Darth Vader is about to roll a **fair** die.
- You have **more uncertainty** about an event because there are **more possibilities**.
- You **receive** more information when you observe it.



$$\begin{aligned} I(x) &= \log_2 \frac{1}{P(6)} \\ &= \log_2 \frac{1}{1/6} \approx \underline{\underline{2.58 \text{ bits}}} \end{aligned}$$

# Information can be additive

- One property of  $I(x) = \log_2 \frac{1}{P(x)}$  is additivity.
- From  $k$  **independent** events  $x_1 \dots x_k$ :
  - Does  $I(x_1 \dots x_k) = I(x_1) + I(x_2) + \dots + I(x_k)$  ?
- The answer is yes!

$$\begin{aligned} I(x_1 \dots x_k) &= \log_2 \frac{1}{P(x_1 \dots x_k)} \\ &= \log_2 \frac{1}{P(x_1) \dots P(x_k)} = \log_2 \frac{1}{P(x_1)} + \dots + \log_2 \frac{1}{P(x_k)} \\ &= I(x_1) + I(x_2) + \dots + I(x_k) \end{aligned}$$

# Aside: Information in computers

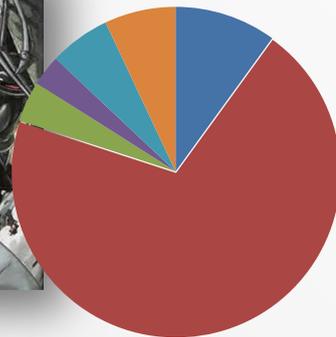
- The unit bit appears familiar to the units describing file sizes...
- And they are related!
- $1\text{ GB} = 2^{10}\text{MB} = 2^{20}\text{KB} = 2^{30}\text{Bytes}$ , where:
  - 1 Byte = 8 bits.
  - Historically: 1 byte was used to store one character.
- File sizes in computers are **described by the amount of information.**
  - The file sizes also depend on the method of encoding (approx. “file format”)

# Events and random variables

- An event  $x$  is a sample from a random variable  $X$ .
- Example 1:
  - $X$ : Darth Vader saying something (either yes or no)
  - $x$ : What DV says ( $x = \text{“no”}$ )
- Example 2:
  - $X$ : Darth Vader rolling a die
  - $x$ : The side facing upwards (e.g.,  $x = 3$ )
- $x$  is deterministic.  $X$  is random.
- $x$  is the output emitted by the “source”  $X$ .

# Information with unequal events

- The random variable  $X$  can take possible values:  $\{v_1, v_2, \dots, v_n\}$ .
- **Each** value has its **own** probability  $\{p_1, p_2, \dots, p_n\}$



■ Yes (0.1)	■ No (0.7)
■ Maybe (0.04)	■ Sure (0.03)
■ Darkside (0.06)	■ Destiny (0.07)

- What is the average amount of information we get in **observing** the **output** of  $X$ ?
- You **still** have 6 events that are possible – **but** you're fairly sure it will be 'No'.

# Entropy

- **Entropy**:  $n.$  the **expected** information gaining from observing the events of the random variable  $X$ .

$$H(X) = E_x[I(x)] = \sum_x p(x) \log \frac{1}{p(x)}$$

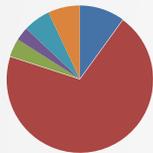
ENTROPY



Notes:

1. Entropy is defined towards a random variable.
2. Entropy is the average uncertainty inherent in a random variable.

# Entropy – examples



- Yes (0.1)
- No (0.7)
- Maybe (0.04)
- Sure (0.03)
- Darkside (0.06)
- Destiny (0.07)

$$H(X) = \sum_i p_i \log_2 \frac{1}{p_i}$$
$$= 0.7 \log_2(1/0.7) + 0.1 \log_2(1/0.1) + \dots$$
$$= 1.542 \text{ bits}$$

There is **less** average uncertainty when the probabilities are ‘skewed’.

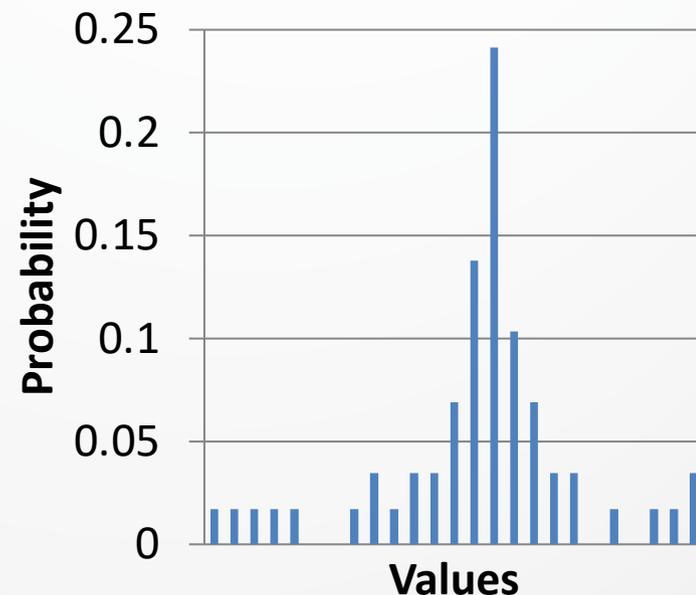
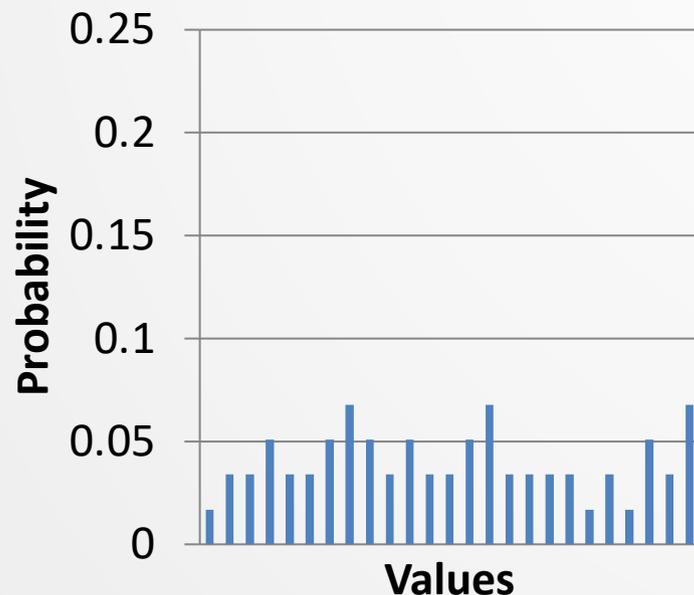


- 1
- 2
- 3
- 4
- 5
- 6

$$H(X) = \sum_i p_i \log_2 \frac{1}{p_i} = 6 \left( \frac{1}{6} \log_2 \frac{1}{1/6} \right)$$
$$= 2.585 \text{ bits}$$

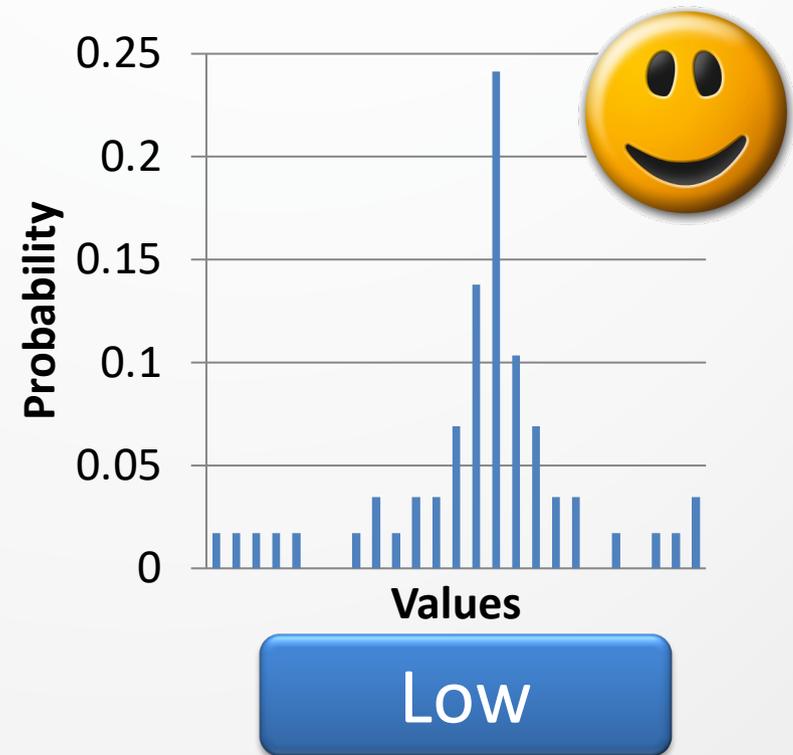
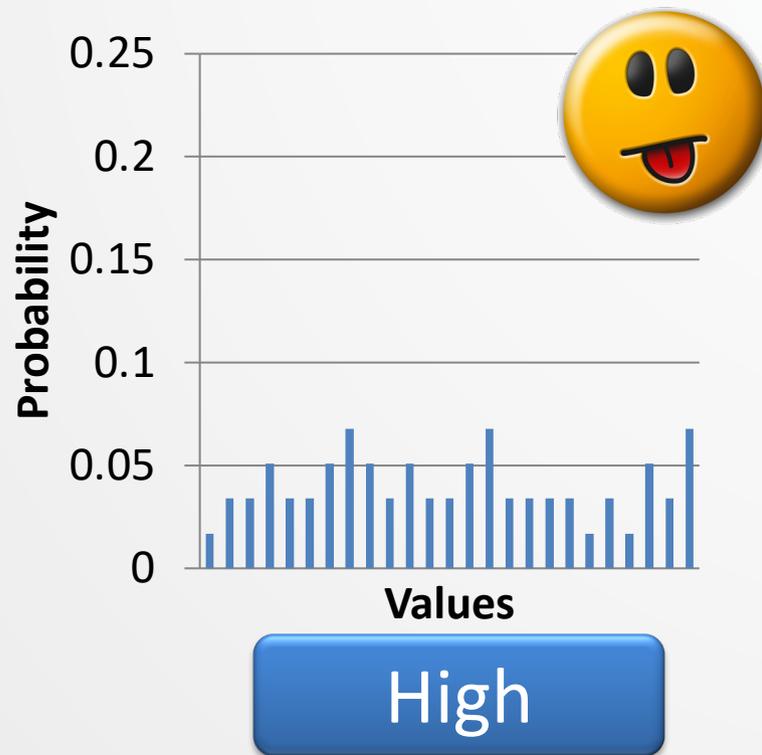
# Entropy characterizes the distribution

- ‘**Flatter**’ distributions have a **higher** entropy because the choices are **more equivalent**, on average.
  - So which of these distributions has a **lower** entropy?



# Low entropy makes decisions easier

- When predicting the next event, we'd like a distribution with **lower** entropy.
  - Low entropy  $\equiv$  less uncertainty

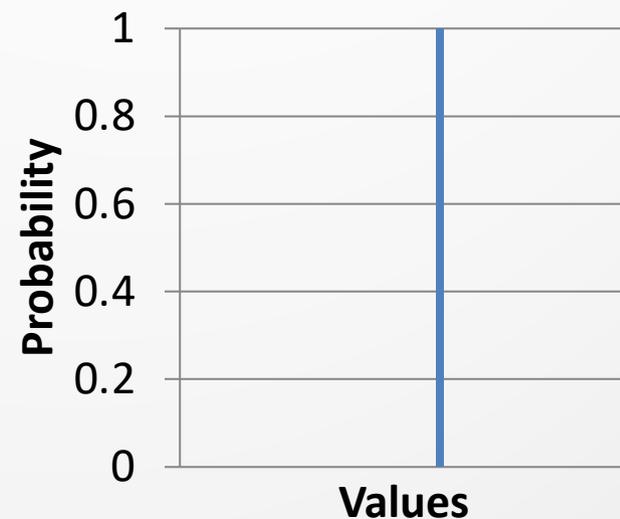
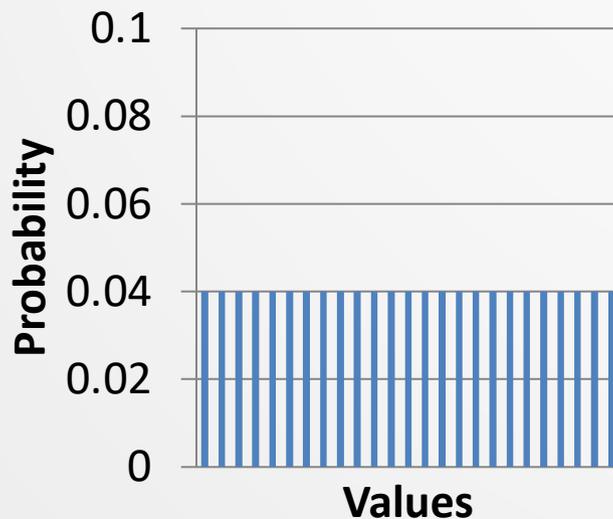


# Bounds on entropy

- **Maximum:** uniform distribution  $X_1$ . Given  $M$  choices,

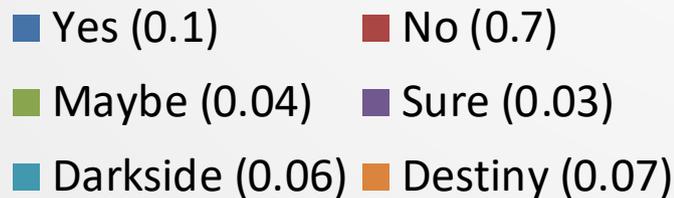
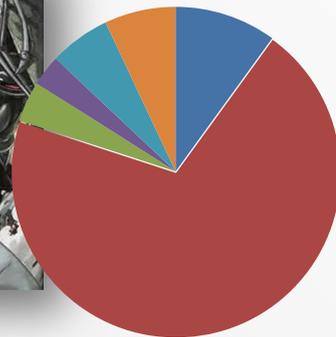
$$H(X_1) = \sum_i p_i \log_2 \frac{1}{p_i} = \sum_i \frac{1}{M} \log_2 \frac{1}{1/M} = \log_2 M$$

- **Minimum:** only one choice,  $H(X_2) = p_i \log_2 \frac{1}{p_i} = 1 \log_2 1 = 0$



# Understanding entropy in coding

- We can **encode** a random variable  $X$  :
  - For a **lossless** encoding,  $X$  can be recovered.
- There are many possible **codes** to encode a random variable.
  - They may involve different **codelengths** (num. bits)

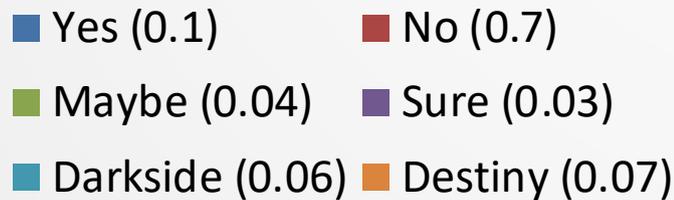
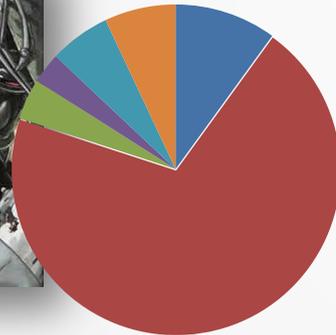


Word (sorted)	Linear Code
No	000
Yes	001
Destiny	010
Darkside	011
Maybe	100
Sure	101

Average codelength = **3 bits**

# Coding with fewer bits is better

- If we want to **transmit** Vader's words **efficiently**, we can **encode** them so that **more probable words** require **fewer bits**.
  - On **average**, fewer bits will need to be transmitted.

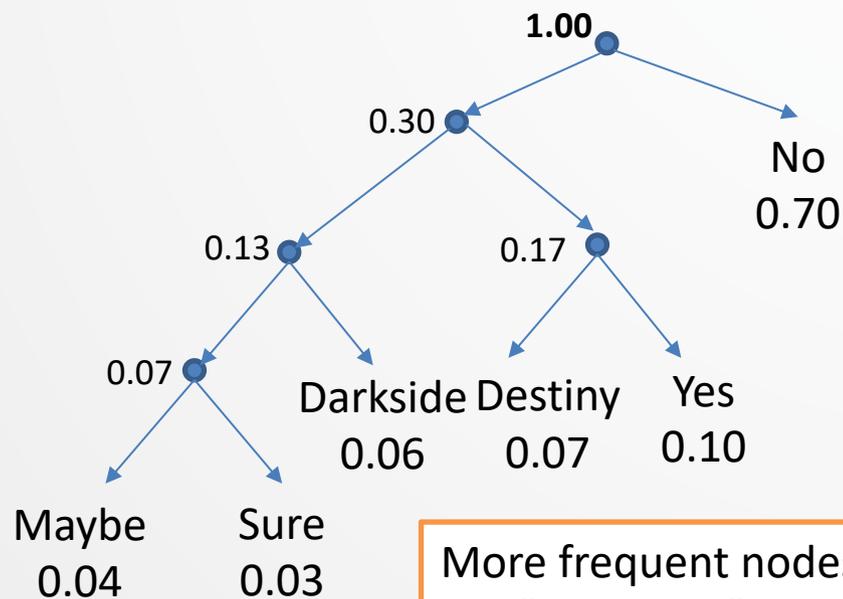


Word (sorted)	Linear Code	Probability	Huffman Code
No	000	0.7	0
Yes	001	0.1	100
Destiny	010	0.07	101
Darkside	011	0.06	110
Maybe	100	0.04	1110
Sure	101	0.03	1111

Average codelength (Huffman) =  $1 * 0.7 + 3 * (0.1 + 0.07 + 0.06) + 4 * (0.04 + 0.03) = 1.67$  bits

# Huffman codes: build tree

- Start with the words: each word is a **leaf node**.
- Merge the two **least** possible nodes into one.
- Repeat until the **Huffman tree** is constructed.



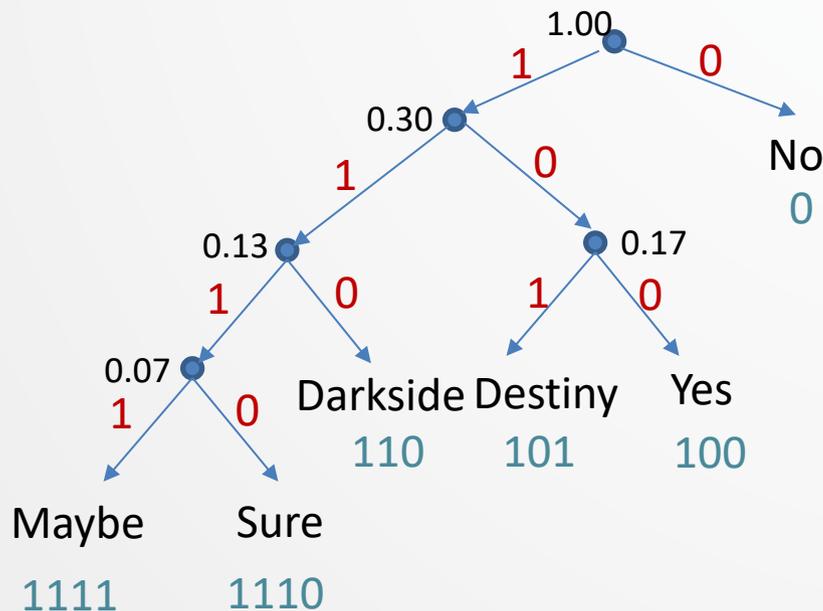
More frequent nodes are “shallower” in the Huffman tree!

Word	Probability
No	0.7
Yes	0.1
Destiny	0.07
Darkside	0.06
Maybe	0.04
Sure	0.03

# Huffman codes: assign codes

- Then assign code values based on the **tree branching**.

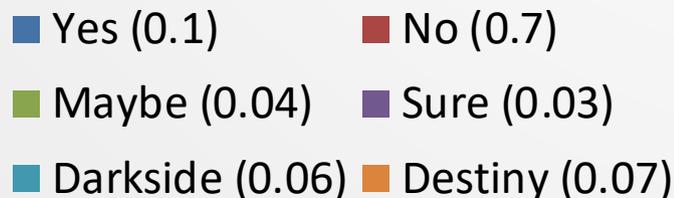
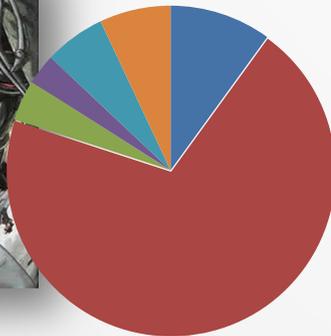
More frequent nodes are assigned shorter codes!



Word	Probability	Huffman Code
No	0.7	0
Yes	0.1	100
Destiny	0.07	101
Darkside	0.06	110
Maybe	0.04	1110
Sure	0.03	1111

# Coding symbols efficiently

- What is the **minimal** possible average **codelength** needed to **losslessly encode** a random variable  $X$ ?
- Answer: entropy!
- $H(X) = \sum_x \log_2 \frac{1}{P(x)} = 1.542 \text{ bits}$



Remark: This is Shannon's  
Source Coding Theorem

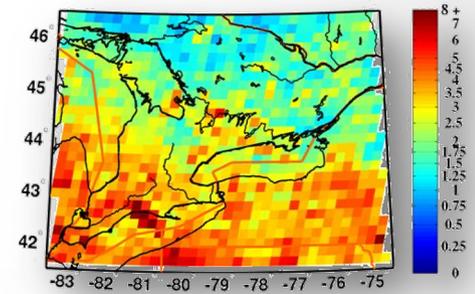
# Alternative notions of entropy

- Entropy is **equivalently**:
  - The **average** amount of **information** provided by an observation of a random variable,
  - The **average** amount of **uncertainty** you have **before** an observation of a random variable,
  - The **average** amount of '**surprise**' you receive during the observation,
  - The number of bits needed to communicate that random variable
    - Aside: Shannon showed that you **cannot** have a **coding scheme** that can communicate it **more efficiently** than  $H(S)$

# Some information-theoretic terms

- Joint entropy
- Conditional entropy
- Mutual information

# Entropy of several variables



- Consider the vocabulary of a meteorologist describing Temperature and Wetness.
  - Temperature = {*hot, mild, cold*}
  - Wetness = {*dry, wet*}

$$P(W = \text{dry}) = 0.6,$$
$$P(W = \text{wet}) = 0.4$$

$$H(W) = 0.6 \log_2 \frac{1}{0.6} + 0.4 \log_2 \frac{1}{0.4} = \mathbf{0.970951 \text{ bits}}$$

$$P(T = \text{hot}) = 0.3,$$
$$P(T = \text{mild}) = 0.5,$$
$$P(T = \text{cold}) = 0.2$$

$$H(T) = 0.3 \log_2 \frac{1}{0.3} + 0.5 \log_2 \frac{1}{0.5} + 0.2 \log_2 \frac{1}{0.2} = \mathbf{1.48548 \text{ bits}}$$

But  $W$  and  $T$  are *not* independent,  
 $P(W, T) \neq P(W)P(T)$

# Joint entropy

- **Joint Entropy:**  $n$ . the **average** amount of information needed to specify **multiple** variables **simultaneously**.

$$H(X, Y) = \sum_x \sum_y p(x, y) \log_2 \frac{1}{p(x, y)}$$

- **Hint:** this is *very* similar to univariate entropy – we just replace univariate probabilities with joint probabilities and sum over everything.

# Entropy of several variables

- Consider joint probability,  $P(W, T)$

	cold	mild	hot	
dry	0.1	0.4	0.1	0.6
wet	0.2	0.1	0.1	0.4
	0.3	0.5	0.2	1.0

- Joint entropy**,  $H(W, T)$ , computed as a sum over the space of joint events ( $W = w, T = t$ )

$$H(W, T) = 0.1 \log_2 1/0.1 + 0.4 \log_2 1/0.4 + 0.1 \log_2 1/0.1 + 0.2 \log_2 1/0.2 + 0.1 \log_2 1/0.1 + 0.1 \log_2 1/0.1 = 2.32193 \text{ bits}$$

Notice  $H(W, T) \approx 2.32 < 2.46 \approx H(W) + H(T)$

# Entropy given knowledge

- In our example, **joint entropy** of two variables together is **lower** than the **sum** of their **individual** entropies
  - $H(W, T) \approx 2.32 < 2.46 \approx H(W) + H(T)$
- **Why?**
- Information is **shared** among variables
  - There are **dependencies**, e.g., between temperature and wetness.
  - E.g., if we knew **exactly** how **wet** it is, is there **less confusion** about what the **temperature** is ... ?

# Conditional entropy

- **Conditional entropy:** *n.* the **average** amount of information needed to specify one variable given that you know another.
  - A.k.a **'equivocation'**

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

- **Hint:** this is *very* similar to how we compute expected values in general distributions.

# Entropy given knowledge

- Consider **conditional** probability,  $P(T|W)$

$P(W, T)$	$T = \text{cold}$	mild	hot	
$W = \text{dry}$	0.1	0.4	0.1	<b>0.6</b>
wet	0.2	0.1	0.1	<b>0.4</b>
	<b>0.3</b>	<b>0.5</b>	<b>0.2</b>	<b>1.0</b>

$$P(T|W) = P(W, T) / P(W)$$

$P(T   W)$	$T = \text{cold}$	mild	hot	
$W = \text{dry}$	0.1/ <b>0.6</b>	0.4/ <b>0.6</b>	0.1/ <b>0.6</b>	<b>1.0</b>
wet	0.2/ <b>0.4</b>	0.1/ <b>0.4</b>	0.1/ <b>0.4</b>	<b>1.0</b>

# Entropy given knowledge

- Consider **conditional** probability,  $P(T|W)$

$P(T W)$	$T = \text{cold}$	mild	hot	
$W = \text{dry}$	1/6	2/3	1/6	1.0
wet	1/2	1/4	1/4	1.0

- $H(T|W = \text{dry}) = H\left(\left\{\frac{1}{6}, \frac{2}{3}, \frac{1}{6}\right\}\right) = 1.25163 \text{ bits}$
- $H(T|W = \text{wet}) = H\left(\left\{\frac{1}{2}, \frac{1}{4}, \frac{1}{4}\right\}\right) = 1.5 \text{ bits}$

- Conditional entropy** combines these:

$$\begin{aligned}
 &H(T|W) \\
 &= [p(W = \text{dry})H(T|W = \text{dry})] + [p(W = \text{wet})H(T|W = \text{wet})] \\
 &= 1.350978 \text{ bits}
 \end{aligned}$$

0.6 0.4

# Equivocation removes uncertainty

- Remember  $H(T) = 1.48548$  bits
  - $H(W, T) = 2.32193$  bits
  - $H(T|W) = 1.350978$  bits
- } Entropy (i.e., confusion) about temperature is **reduced** if we **know** how wet it is outside.
- How much does  $W$  tell us about  $T$ ?
    - $H(T) - H(T|W) = 1.48548 - 1.350978 \approx 0.1345$  bits
    - Well, a little bit!

# Perhaps $T$ is more informative?

- Consider **another** conditional probability,  $P(W|T)$

$P(W T)$	$T = \text{cold}$	mild	hot
$W = \text{dry}$	0.1/ <b>0.3</b>	0.4/ <b>0.5</b>	0.1/ <b>0.2</b>
wet	0.2/ <b>0.3</b>	0.1/ <b>0.5</b>	0.1/ <b>0.2</b>
	1.0	1.0	1.0

- $H(W|T = \text{cold}) = H\left(\left\{\frac{1}{3}, \frac{2}{3}\right\}\right) = 0.918295$  bits
- $H(W|T = \text{mild}) = H\left(\left\{\frac{4}{5}, \frac{1}{5}\right\}\right) = 0.721928$  bits
- $H(W|T = \text{hot}) = H\left(\left\{\frac{1}{2}, \frac{1}{2}\right\}\right) = 1$  bit
- $H(W|T) = 0.8364528$  bits**

# Equivocation removes uncertainty

- $H(T) = 1.48548$  bits
- $H(W) = 0.970951$  bits
- $H(W, T) = 2.32193$  bits
- $H(T|W) = 1.350978$  bits
- $H(T) - H(T|W) \approx \mathbf{0.1345 \text{ bits}}$   Previously computed

- How much does  $T$  tell us about  $W$  on average?
  - $H(W) - H(W|T) = 0.970951 - 0.8364528$   
 $\approx \mathbf{0.1345 \text{ bits}}$
- Interesting ... is that a coincidence?

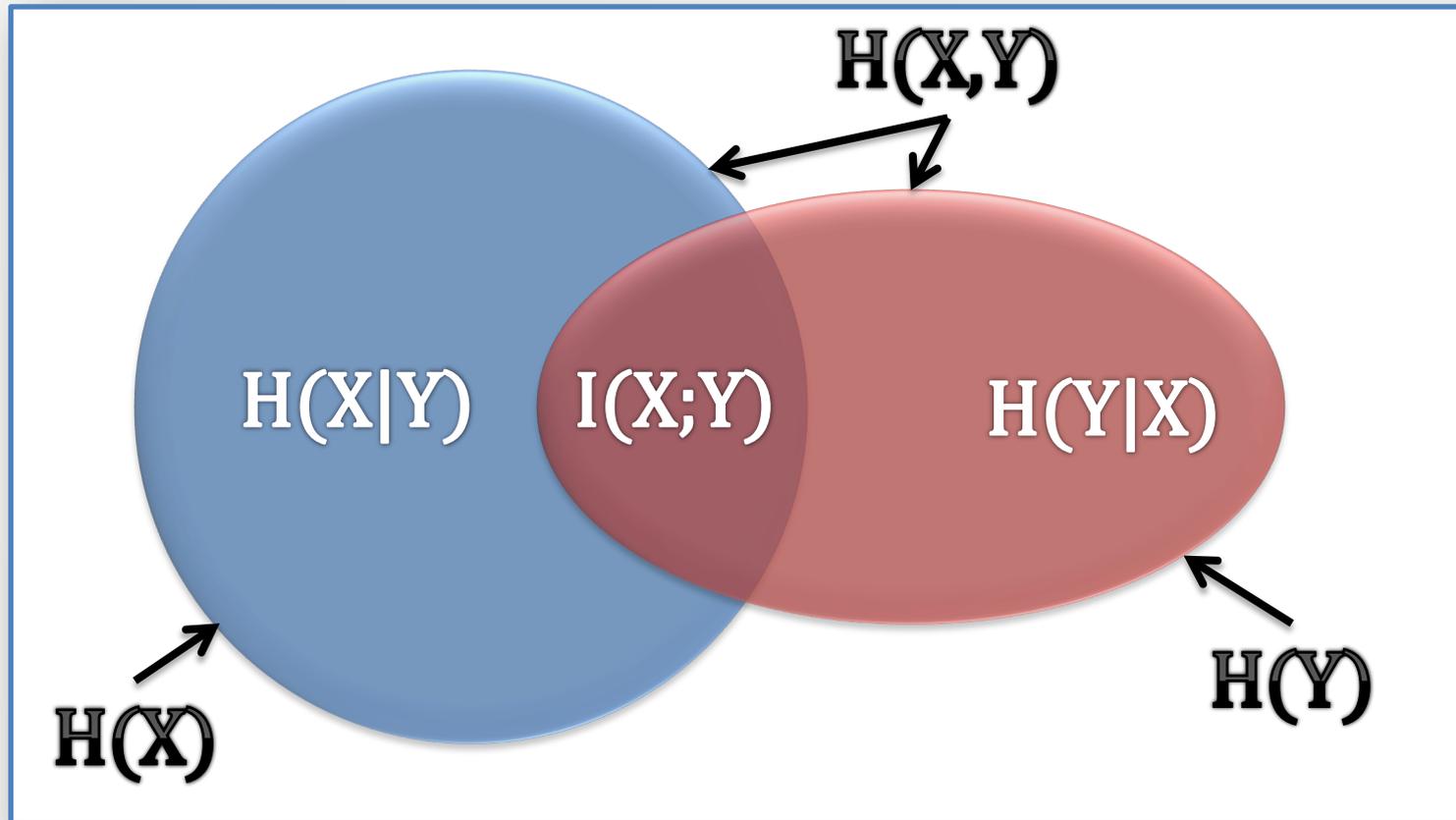
# Mutual information

- **Mutual information:**  $n$ . the **average** amount of information **shared** between variables.

$$\begin{aligned} I(X; Y) &= H(X) - H(X|Y) = H(Y) - H(Y|X) \\ &= \sum_{x,y} p(x, y) \log_2 \frac{p(x, y)}{p(x)p(y)} \end{aligned}$$

- **Hint:** The amount of uncertainty **removed** in variable  $X$  if you know  $Y$ .
- **Hint2:** If  $X$  and  $Y$  are **independent**,  $p(x, y) = p(x)p(y)$ , then
$$\log_2 \frac{p(x, y)}{p(x)p(y)} = \log_2 1 = 0 \quad \forall x, y - \text{there is no mutual information!}$$

# Relations between entropies

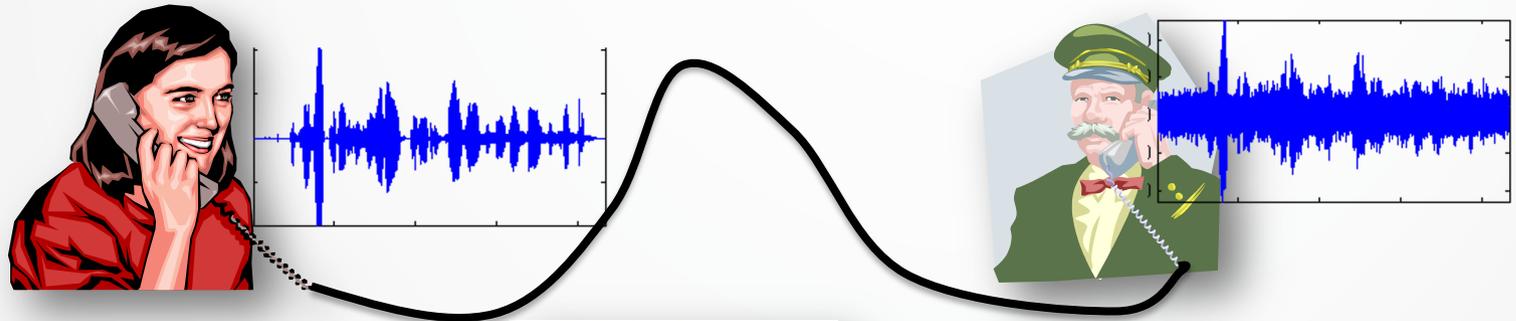


$$H(X, Y) = H(X) + H(Y) - I(X; Y)$$

# Reminder – the noisy channel

- Messages can get **distorted** when passed through a **noisy** conduit – how much information is lost/retained?

- Signals



- Symbols



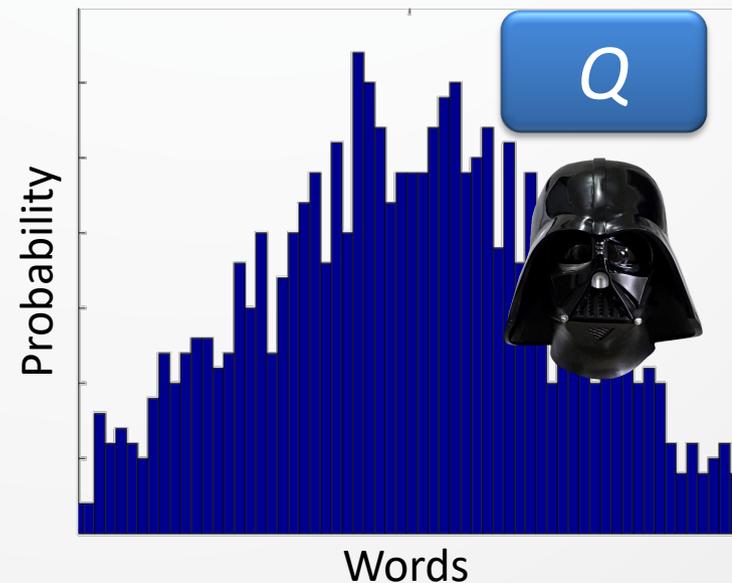
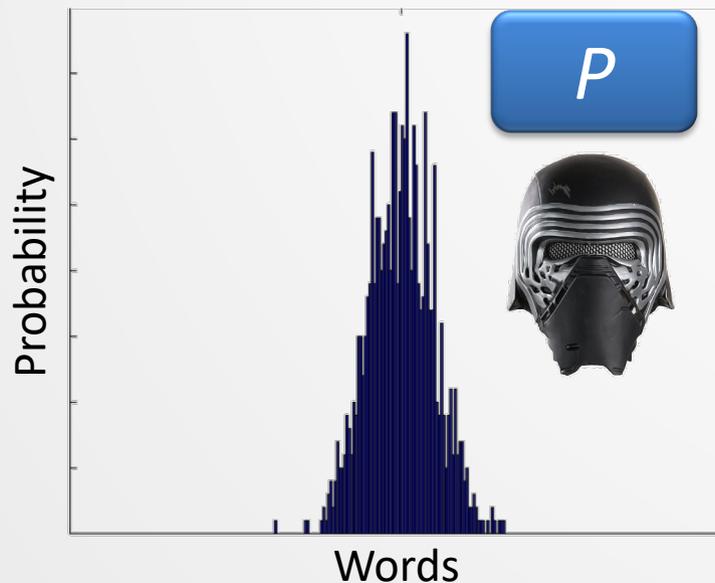
- Languages



# Relating corpora

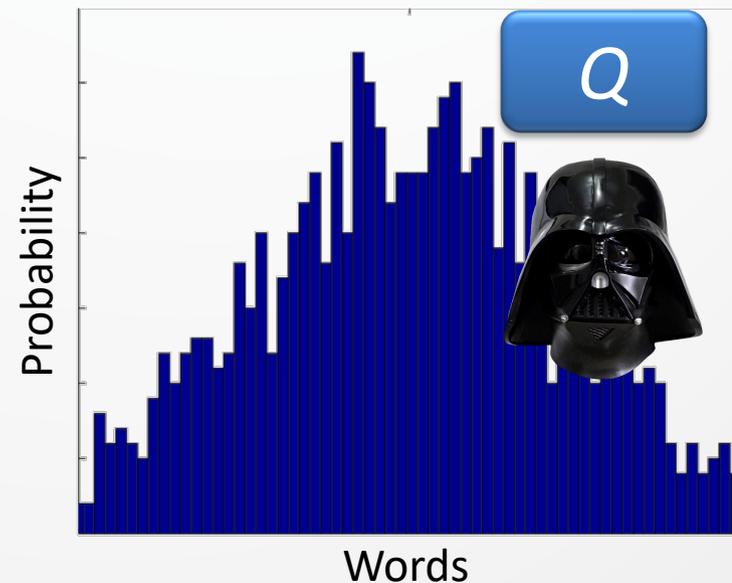
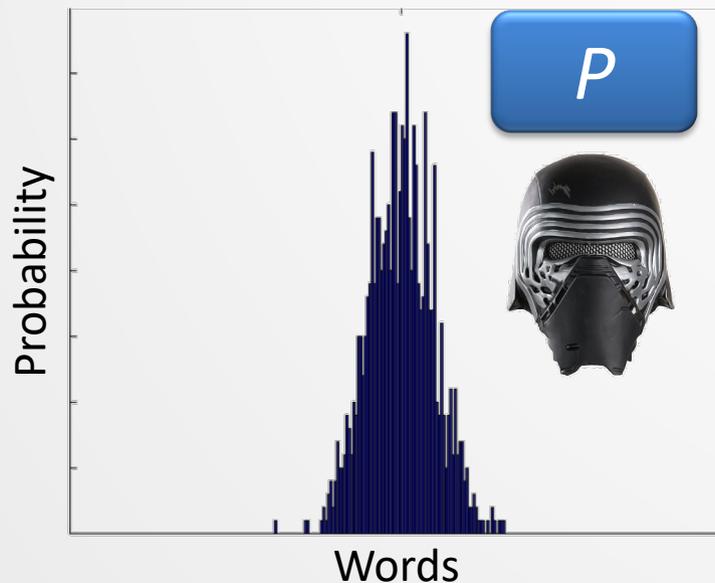
# Relatedness of two distributions

- How **similar** are two probability distributions?
  - e.g., Distribution  $P$  learned from *Kylo Ren*  
Distribution  $Q$  learned from *Darth Vader*



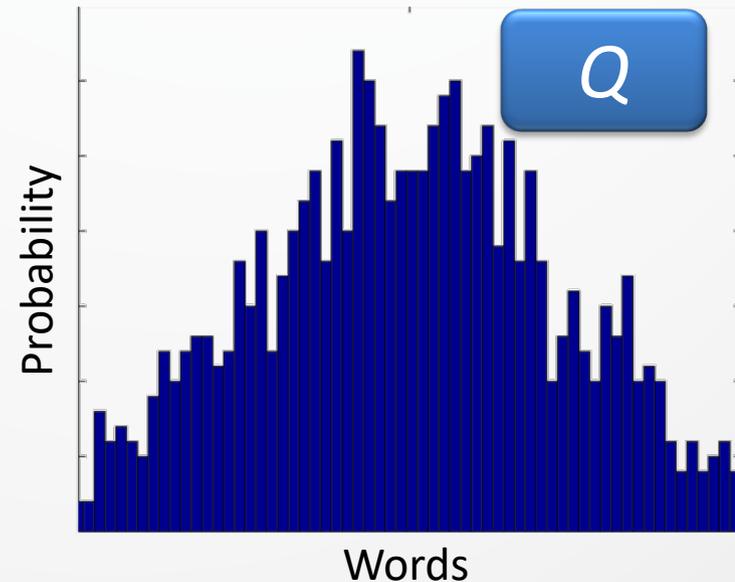
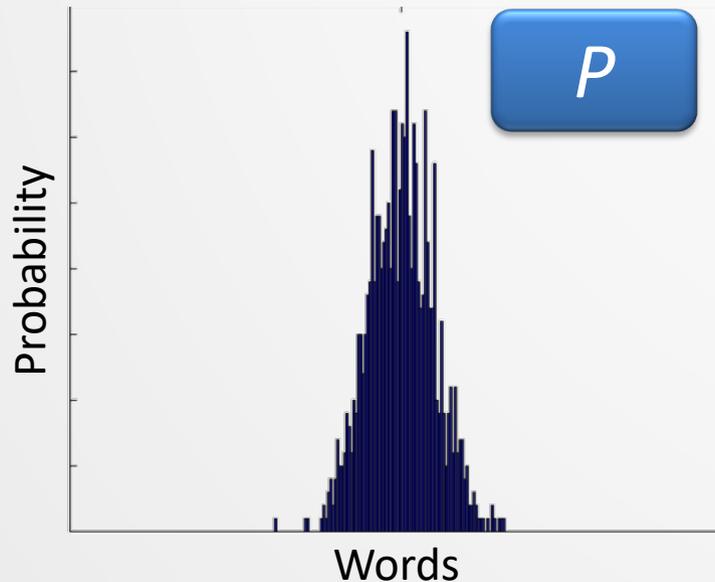
# Relatedness of two distributions

- A Huffman code based on Vader ( $Q$ ) instead of Kylo ( $P$ ) will be less *efficient* at coding symbols that Kylo will say.
- What is the **average number of extra bits** required to code symbols from  $P$  when using a code based on  $Q$ ?



# Kullback-Leibler divergence

- **KL divergence:**  $n$ . the **average log difference** between the distributions  $P$  and  $Q$ , relative to  $Q$ .  
a.k.a. **relative entropy**.  
*caveat:* we assume  $0 \log 0 = 0$



# Kullback-Leibler divergence

$$D_{KL}(P||Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

- Why  $\log \frac{P(i)}{Q(i)}$ ?
- $\log \frac{P(i)}{Q(i)} = \log P(i) - \log Q(i) = \log \left( \frac{1}{Q(i)} \right) - \log \left( \frac{1}{P(i)} \right)$
- If word  $w_i$  is less probable in  $Q$  than  $P$  (i.e., it carries more information), it will be Huffman encoded in more bits, so when we see  $w_i$  from  $P$ , we need  $\log \frac{P(i)}{Q(i)}$  more bits.

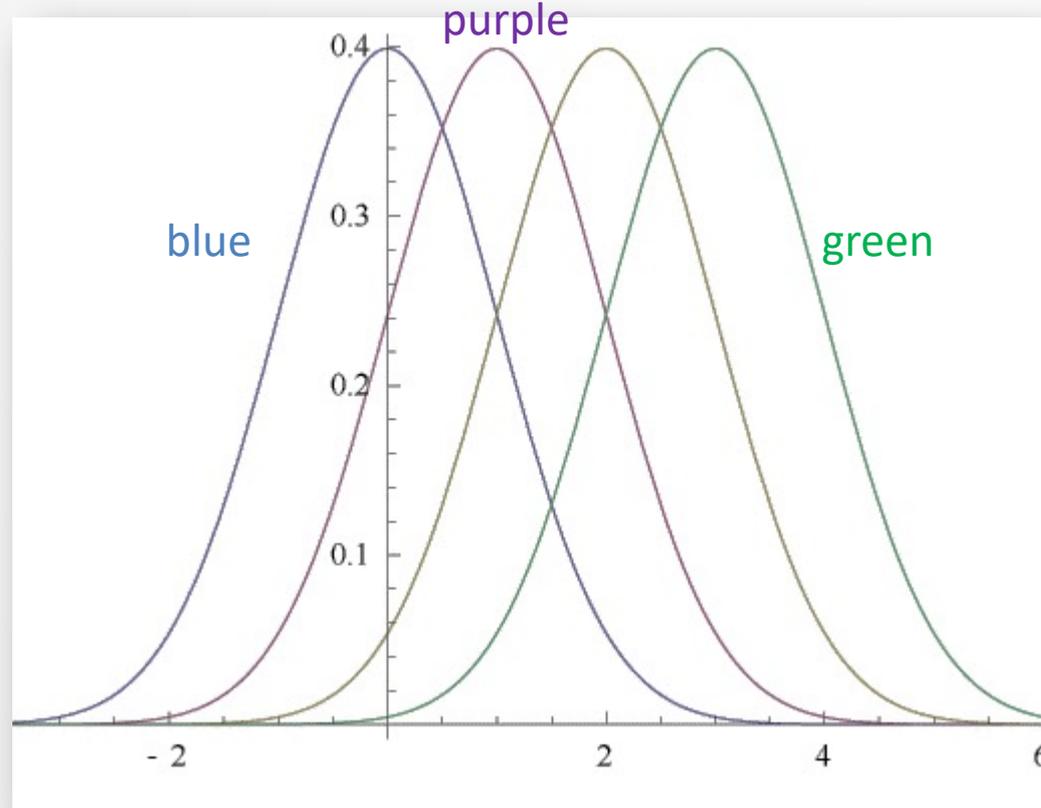
# Kullback-Leibler divergence

- KL divergence:
  - is *somewhat* like a '**distance**' :
    - $D_{KL}(P||Q) \geq 0 \quad \forall P, Q$
    - $D_{KL}(P||Q) = 0$  iff  $P$  and  $Q$  are identical.
  - is **not symmetric**,  $D_{KL}(P||Q) \neq D_{KL}(Q||P)$
- Aside 1: Jensen-Shannon divergence is symmetric.
- Aside 2:

$$I(P; Q) = D_{KL}(P(X, Y)||P(X)P(Y))$$

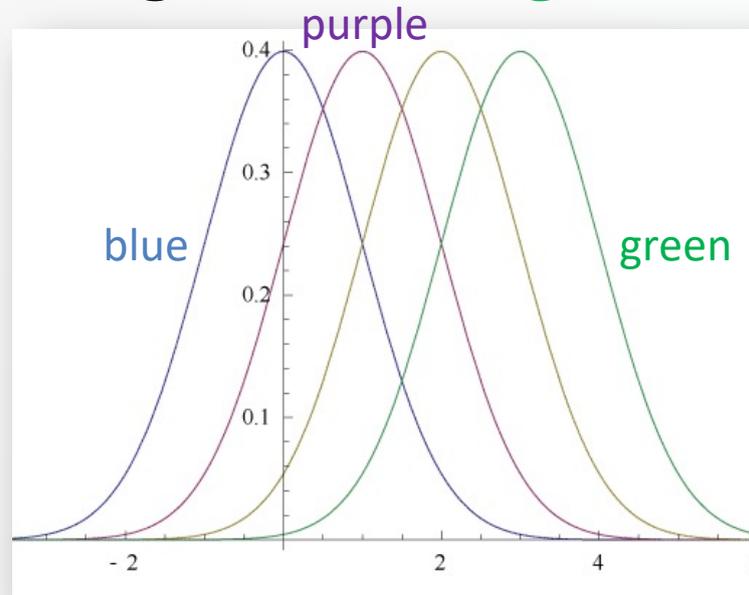
# Kullback-Leibler divergence

- KL divergence generalizes to **continuous** distributions.
- Below,  $D_{KL}(\textit{blue}||\textit{green}) > D_{KL}(\textit{blue}||\textit{purple})$



# Applications of KL divergence

- Often used towards some **other purpose**, e.g.,
  - In **evaluation** to say that *purple* is a **better** model than *green* of the **true distribution** *blue*.
  - In **machine learning** to adjust the parameters of *purple* to be, e.g., less like *green* and more like *blue*.



# Entropy as intrinsic LM evaluation

- **Cross-entropy** measures how difficult it is to encode an event drawn from a **true probability**  $p$  given a **model** based on a distribution  $q$ .

- What if we don't know the **true probability**  $p$ ?
  - We'd have to estimate the CE using a test corpus  $C$ :

$$H(p, q) \approx - \frac{\log_2 P_q(C)}{\|C\|}$$

- What's the probability of a corpus  $P_q(C)$ ?

# Probability of a corpus?

- The probability  $P(C)$  of a **corpus**  $C$  requires similar **assumptions** that allowed us to compute the probability  $P(s_i)$  of a **sentence**  $s_i$ .

	Sentence	Corpus
Chain rule	$P(s_i) = P(w_1) \prod_{t=2}^n P(w_t   w_{1:(t-1)})$	$P(C) = P(w_1) \prod_{t=2}^{\ C\ } P(w_t   w_{1:(t-1)})$
Approx.	$P(s_i) \approx \prod_t P(w_t)$	$P(C) \approx \prod_i P(s_i)$

- Regardless** of the LM used for  $P(s_i)$ , we can assume **complete independence** between sentences.

# Intrinsic evaluation – Cross-entropy

- **Cross-entropy** of a LM  $M$  and a *new* test corpus  $C$  with size  $\|C\|$  (total number of words), where sentence  $s_i \in C$ , is *approximated* by:

$$H(C; M) = -\frac{\log_2 P_M(C)}{\|C\|} = -\frac{\sum_i \log_2 P_M(s_i)}{\sum_i \|s_i\|}$$

- **Perplexity** comes from this definition:

$$PP_M(C) = 2^{H(C; M)}$$

# Cross-entropy in Machine Learning

- **Cross-entropy** in ML measures the quality of a predicted distribution  $q(Y)$  with respect to  $p(Y)$ :

$$H(p, q) = \sum_y p(y) \log \frac{1}{q(y)}$$

- Note 1: ML usually uses log with base e.
- Note 2: Cross entropy is usually used as the target for optimization, i.e., **cross-entropy loss**.
- Note 3: This is also called **log-loss**, or **negative log-likelihood loss**.

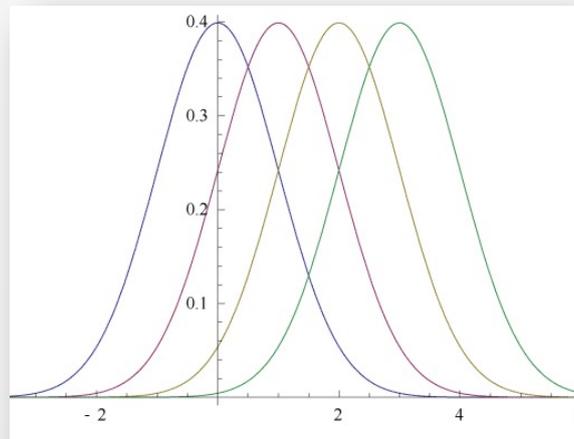
# Decisions

# Deciding what we know

- **Anecdotes** are often useless except as proofs by contradiction.
  - E.g., “*I saw Google used as a verb*” does **not** mean that *Google* is **always** (or even **likely** to be) a verb, just that it is **not always** a noun.
- **Shallow statistics** are often not enough to be truly meaningful.
  - E.g., “*My ASR system is 95% accurate on my test data. Yours is only 94.5% accurate, you horrible knuckle-dragging idiot.*”
    - What if the test data was **biased** to favor my system?
    - What if we only used a **very small** amount of data?
- Given all this potential ambiguity, we need a **test** to see if our statistics actually **mean** something.

# Differences due to sampling

- We saw that **KL divergence** essentially measures how **different** two distributions are from each other.
- But what if their difference is due to **randomness in sampling**?
- How can we tell that a distribution is **really** different from another?



# Hypothesis testing

- Often, we assume a **null hypothesis**,  $H_0$ , which states that the **two distributions are the same** (i.e., come from the same underlying model, population, or phenomenon).
- We **reject** the null hypothesis if the probability of it being true is too small.
  - This is often our goal – e.g., if my ASR system beats yours by 0.5%, I want to show that this difference is **not** a random accident.
  - I assume it *was* an accident, then show how nearly *impossible* that is.
  - As scientists, we have to be very **careful** to not reject  $H_0$  too hastily.
    - How can we ensure our **diligence**?

# Confidence

- We **reject**  $H_0$  if it is **too improbable** based on the evidence.
  - How do we determine the value of ‘too’?
- **Significance level**  $\alpha$  ( $0 \leq \alpha \leq 1$ ) is the **maximum** probability that two distributions are **identical** allowing us to **disregard**  $H_0$ .
  - In practice,  $\alpha \leq 0.05$ . Usually, it’s much lower.
  - **Confidence level** is  $\gamma = 1 - \alpha$
  - E.g., a confidence level of **95%** ( $\alpha = 0.05$ ) implies that we expect that our decision is correct 95% of the time, **regardless of the test data.**

# Confidence

- We will briefly see three types of **statistical tests** that can tell us how **confident** we can be in a claim:
  1. A **t-test**, which usually tests whether the **means** of two models are the same. There are many types, but most assume **Gaussian** distributions.
  2. An **analysis of variance (ANOVA)**, which generalizes the *t*-test to more than two groups.
  3. The  **$\chi^2$  test**, which evaluates **categorical** (discrete) outputs.

# 1. The $t$ -test

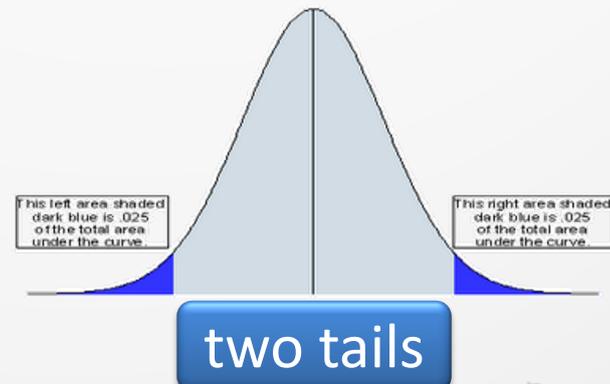
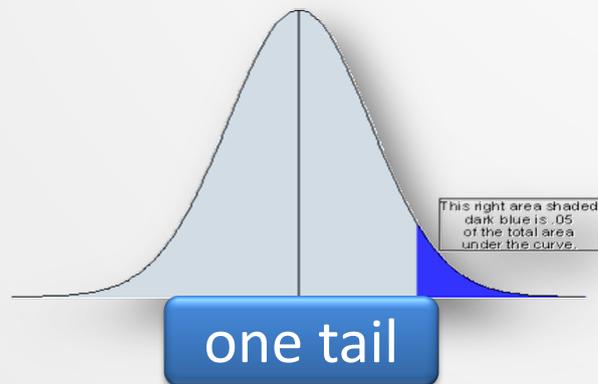
- The  $t$ -test is a method to compute if distributions are significantly different from one another.
- It is based on the mean ( $\bar{x}$ ) and variance ( $\sigma^2$ ) of  $N$  samples.
- It compares  $\bar{x}$  and  $\sigma$  to  $H_0$  which states that the samples are drawn from a distribution with a mean  $\mu$ .
- If  $t = \frac{\bar{x} - \mu}{\sqrt{\sigma^2 / N}}$  (the “t-statistic”) is large enough, we can reject  $H_0$ .

An example would be nice...

There are actually **several types** of  $t$ -tests for different situations...

# Example of the $t$ -test: tails

- Imagine the average tweet length of a McGill 'student' is  $\mu = 158$  chars.
- We sample  $N = 200$  UofT students and find that our average tweet is  $\bar{x} = 169$  chars (with  $\sigma^2 = 2600$ ).
- Are UofT tweets significantly **longer** than McGill tweets?
- We use a '**one-tailed**' test because we want to see if UofT tweet lengths are significantly **higher**.
  - If we just wanted to see if UofT tweets were significantly **different**, we'd use a **two-tailed** test.



# Example of the $t$ -test: freedom

- Imagine the average tweet length of a McGill ‘student’ is  $\mu = 158$  chars.
- We sample  $N = 200$  UofT students and find that our average tweet is  $\bar{x} = 169$  chars (with  $\sigma^2 = 2600$ ).
- Are UofT tweets significantly **longer** than McGill tweets?
- **Degrees of freedom (d.f.):** *n.pl.* In *this*  $t$ -test, this is the sum of the number of observations, minus 1 (the number of sample sets).
- In our example, we have  $N_{UofT} = 200$  for UofT students, meaning  
 $d.f. = 199$ 
  - (this example is adapted from Manning & Schütze)

# Example of the $t$ -test

- Imagine the average tweet length of a McGill 'student' is  $\mu = 158$  chars.
- We sample  $N = 200$  UofT students and find that our average tweet is  $\bar{x} = 169$  chars (with  $\sigma^2 = 2600$ ).
- Are UofT tweets significantly **longer** than McGill tweets?
- So  $t = \frac{\bar{x} - \mu}{\sqrt{\sigma^2 / N}} = \frac{169 - 158}{\sqrt{2600 / 200}} \approx 3.05$
- In a  **$t$ -test table**, we look up the minimum value of  $t$  necessary to reject  $H_0$  at  $\alpha = 0.005$  (we want to be quite confident) for a 1-tailed test...

# Example of the $t$ -test

- So  $t = \frac{\bar{x} - \mu}{\sqrt{\sigma^2/N}} = \frac{169 - 158}{\sqrt{2600/200}} \approx 3.05$
- In a  **$t$ -test table**, we look up the minimum value of  $t$  necessary to reject  $H_0$  at  $\alpha = 0.005$ , and find 2.576 (using  $d.f. = 199 \approx \infty$ )
  - Since  $3.05 > 2.576$ , we can reject  $H_0$  at the 99.5% level of confidence ( $\gamma = 1 - \alpha = 0.995$ ); **UofT students are significantly more verbose.**

	$\alpha$ (one-tail)	0.05	0.025	0.01	0.005	0.001	0.0005
d.f.	1	6.314	12.71	31.82	63.66	318.3	636.6
	10	1.812	2.228	2.764	3.169	4.144	4.587
	20	1.725	2.086	2.528	2.845	3.552	3.850
	$\infty$	1.645	1.960	2.326	<b>2.576</b>	3.091	3.291

# Example of the $t$ -test

- Some things to observe about the  $t$ -test table:
  - We need **more evidence,  $t$** , if we want to be **more confident** (left-right dimension).
  - We need **more evidence,  $t$** , if we have **fewer measurements** (top-down dimension).
- A common criticism of the  $t$ -test is that picking  $\alpha$  is ad-hoc. There are ways to correct for the selection of  $\alpha$ .

	$\alpha$ (one-tail)	0.05	0.025	0.01	0.005	0.001	0.0005
d.f.	1	6.314	12.71	31.82	63.66	318.3	636.6
	10	1.812	2.228	2.764	3.169	4.144	4.587
	20	1.725	2.086	2.528	2.845	3.552	3.850
	$\infty$	1.645	1.960	2.326	<b>2.576</b>	3.091	3.291

# Another example: collocations

- **Collocation:** *n.* a ‘turn-of-phrase’ or usage where a sequence of words is ‘**perceived**’ to have a meaning ‘**beyond**’ the sum of its parts.
- E.g., ‘*disk drive*’, ‘*video recorder*’, and ‘*soft drink*’ are collocations. ‘*cylinder drive*’, ‘*video storer*’, ‘*weak drink*’ are **not** despite some near-synonymy between alternatives.
- Collocations are **not** just highly frequent bigrams, otherwise ‘*of the*’, and ‘*and the*’ would be collocations.
- How can we test if a bigram is a collocation or not?

# Hypothesis testing collocations

- For collocations, the **null hypothesis**  $H_0$  is that there is **no association** between two given words **beyond pure chance**.
  - I.e., the bigram's **actual** distribution and pure chance are the **same**.
  - We compute the probability of those words occurring together if  $H_0$  were true. If that probability is **too low**, we **reject**  $H_0$ .
- E.g., we expect '*of the*' to occur together, because they're both likely words to draw randomly
  - We could probably **not** reject  $H_0$  in that case.

# Example of the *t*-test on collocations

- Is '*new companies*' a collocation?
- In our corpus of 14,307,668 word tokens, *new* appears 15,828 times and *companies* appears 4,675 times.
- Our **null hypothesis**,  $H_0$  is that they are **independent**, i.e.,

$$\begin{aligned} H_0: P(\text{new companies}) &= P(\text{new})P(\text{companies}) \\ &= \frac{15828}{14307668} \times \frac{4675}{14307668} \\ &\approx 3.615 \times 10^{-7} \end{aligned}$$

# Example of the $t$ -test on collocations

- The Manning & Schütze text claims that if the process of randomly generating bigrams follows a **Bernoulli distribution**.
  - i.e., assigning 1 whenever *new companies* appears and 0 otherwise gives  $\bar{x} = p = P(\text{new companies})$
  - For Bernoulli distributions,  $\sigma^2 = p(1 - p)$ . Manning & Schütze claim that we can assume  $\sigma^2 = p(1 - p) \approx p$ , since for most bigrams,  $p$  is very small.

# Example of the $t$ -test on collocations

- So,  $\mu = 3.615 \times 10^{-7}$  is the expected mean in  $H_0$ .
- We **actually count** 8 occurrences of *new companies* in our corpus
  - $\bar{x} = \frac{8}{14307667} \approx 5.591 \times 10^{-7}$ 

There is 1 fewer bigram instance than word tokens in the corpus

$\therefore \sigma^2 \approx p = \bar{x} = 5.591 \times 10^{-7}$
- So  $t = \frac{\bar{x} - \mu}{\sqrt{\sigma^2 / N}} = \frac{5.591 \times 10^{-7} - 3.615 \times 10^{-7}}{\sqrt{5.591 \times 10^{-7} / 14307667}} \approx 0.9999$
- In a  **$t$ -test table**, we look up the minimum value of  $t$  necessary to reject  $H_0$  at  $\alpha = 0.005$ , and find **2.576**.
  - Since **0.9999** < **2.576**, we cannot reject  $H_0$  at the 99.5% level of confidence.
    - We **don't have enough evidence** to think that *new companies* is a collocation (we can't say that it definitely *isn't*, though!).

# Types of $t$ -tests

- We usually use three types of  $t$ -tests:
- **One-sample  $t$ -test**: whether a variable  $X$  equals a known value  $\mu$ .
  - Both the previous two examples are one-sample  $t$ -tests.
  - $X$  is a **random variable**: e.g., mean Tweet length of UofT students.
  - $\mu$  is a specified **constant**. E.g., 0.
- **Two-sample  $t$ -test**: whether a variable  $X$  equals another variable  $Y$ .
  - Example:  $X/Y$ : mean of UofT/McGill tweet lengths.
  - Two-sample  $t$ -test is useful when you sample from **both** UofT and McGill students.
- **Paired  $t$ -test**: whether  $X - Y$  equals a known value  $\mu$ .
  - Example:  $X/Y$ : weight of the participant before/after an exercise. Test whether the exercise reduces weight.
  - Paired  $t$ -test is just one-sample  $t$ -test on the difference (i.e.,  $X - Y$ )
  - Paired  $t$ -test is useful when individual effects matter.

# The normality assumption of $t$ -test

- $t$ -tests assumes the **random variables** are **normally distributed**.
  - Without a normality assumption, don't use  $t$ -tests...
  - You can use **nonparametric** tests instead, e.g., **Mann-Whitney U test**
  - Next slide: For other tests, we can use ANOVA and  $\chi^2$  tests too.
- Usually, the normality is supported by the **central limit theorem**:
  - The **mean** of  $n \rightarrow \infty$  **independent** samples from **any** distribution approximates a normal distribution (details omitted)
- There are some tests to check normality.
  - E.g., **Shapiro-Wilks test**
  - As an exploratory analysis, just do a **quantile-quantile plot (Q-Q plot)** against a normal distribution.

Most tests are one-liners in either **scipy** or **scikit-learn**

## 2. Analysis of variance (aside)

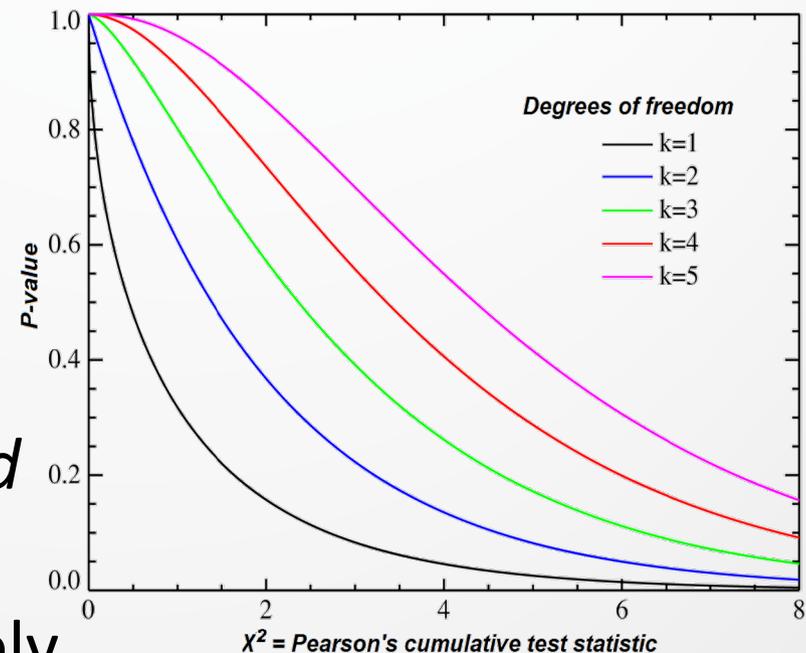
- **Analyses of variance (ANOVAs)** (there are several types) can be:
  - A way to **generalize *t*-tests** to more than two groups.
  - A way to **determine which** (if any) of several **variables** are **responsible** for the **variation** in an observation (and the interaction between them).
- An ANOVA usually involves these steps:
  - Compute a statistic,  $F$ .
  - Which is, *approximately*, a ratio between two variances (divided by their degrees of freedom).
  - The  $F$  statistic, together with the two degrees of freedom, gives us a  $p$  value.
  - If this  $p$  value is smaller than the significance level  $\alpha$ , reject  $H_0$ .

# 3. Pearson's $\chi^2$ test

- The  $\chi^2$  test applies to **categorical** data, like the output of a **classifier**.
- Like the  $t$ -test, we decide on the degrees of freedom (number of categories minus number of parameters), compute the test-statistic, then look it up in a table.
- The test statistic is:

$$\chi^2 = \sum_{c=1}^C \frac{(O_c - E_c)^2}{E_c}$$

where  $O_c$  and  $E_c$  are the *observed* and *expected* number of observations of type  $c$ , respectively.



# 3. Pearson's $\chi^2$ test



- For example, is the die of Darth Vader fair or not?
- Imagine we throw it 60 times. The expected number of appearances of each side is 10.

$c$	$O_c$	$E_c$	$O_c - E_c$	$(O_c - E_c)^2$	$(O_c - E_c)^2 / E_c$
1	5	10	-5	25	2.5
2	8	10	-2	4	0.4
3	9	10	-1	1	0.1
4	8	10	-2	4	0.4
5	10	10	0	0	0
6	20	10	10	100	10
Sum ( $\chi^2$ )					<b>13.4</b>

- With  $df = 6 - 1 = 5$ , the critical value is  $11.07 < \mathbf{13.4}$ , so we throw away  $H_0$ : the die is biased.
- We'll see  $\chi^2$  again soon...

# Entropy and decisions

- **Information theory** is a vast ocean that provides statistical models of communication at the heart of **cybernetics**.
  - We've only taken a first step on the beach.
  - See the ground-breaking work of Shannon & Weaver, e.g.
- So far, we've mainly dealt with **random variables** that the world provides – e.g., words tokens, mainly.
- What if we could transform those inputs into new random variables, or **features**, that are directly engineered to be useful to decision tasks...

# Feature selection

# Determining a good set of features

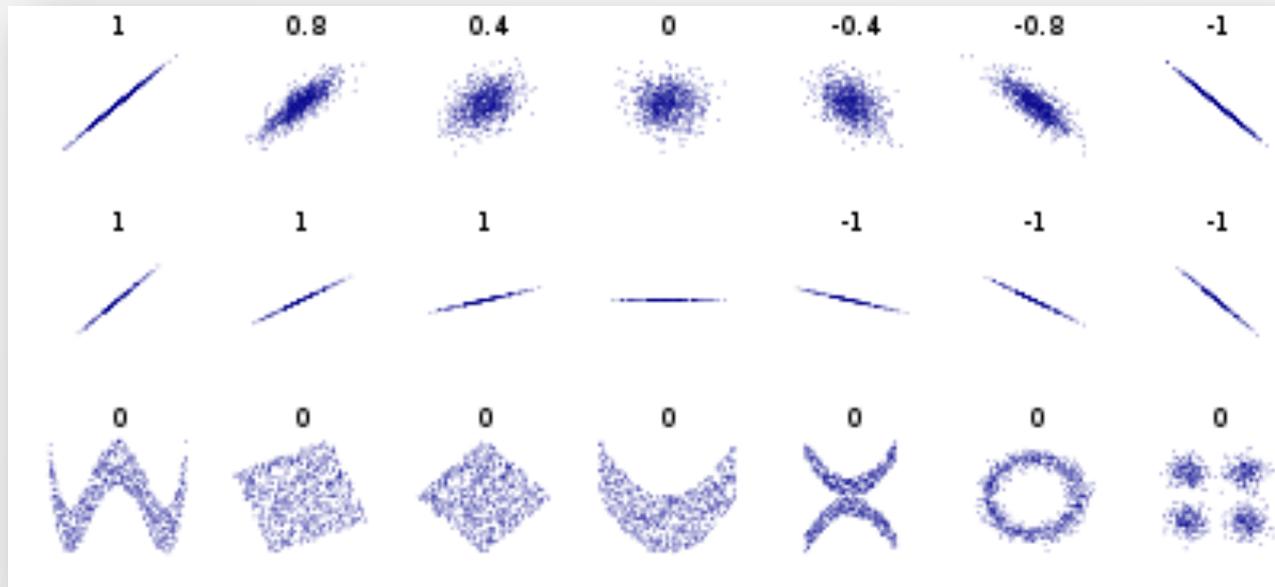
- **Restricting** your feature set to a proper subset quickens **training** and reduces **overfitting**.
- There are a few methods that select good features, e.g.,
  1. Correlation-based feature selection
  2. Minimum Redundancy, Maximum Relevance
  3.  $\chi^2$

# 1. Pearson's correlation

- **Pearson** is a measure of **linear** dependence

$$\rho_{XY} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

- Does not measure '**slope**' nor **non-linear** relations.



# 1. Spearman's correlation

- **Spearman** is a non-parametric measure of **rank** correlation,  $r_{cX} = r(c, X)$ .
  - It is basically Pearson's correlation, but on 'rank variables' that are monotonically increasing integers.
  - If the class  $c$  can be **ordered** (e.g., in any binary case), then we can compute the correlation between a feature  $X$  and that class.

# 1. Correlation-based feature selection

- ‘Good’ features should correlate **strongly** (+ or -) with the *predicted variable* but **not** with other *features*.
- $S_{CFS}$  is some set  $S$  of  $k$  features  $f_i$  that maximizes this ratio, given class  $c$ :

$$S_{CFS} = \operatorname{argmax}_S \frac{\sum_{f_i \in S} r_{cf_i}}{\sqrt{k + 2 \sum_{i=1}^{k-1} \sum_{j=i+1}^k \rho_{f_i f_j}}}$$

## 2. mRMR feature selection

- **Minimum-redundancy-maximum-relevance (mRMR)** can use **correlation**, **distance** scores (e.g.,  $D_{KL}$ ) or **mutual information** to select features.
- For feature set  $S$  of features  $f_i$ , and class  $c$ ,  
 $D(S, c)$  : a measure of **relevance**  $S$  has for  $c$ , and  
 $R(S)$  : a measure of the **redundancy** within  $S$ ,

$$S_{mRMR} = \operatorname{argmax}_S [D(S, c) - R(S)]$$

## 2. mRMR feature selection

- Measures of **relevance** and **redundancy** can make use of our familiar measures of *mutual information*,

- $$D(S, c) = \frac{1}{\|S\|} \sum_{f_i \in S} I(f_i; c)$$

- $$R(S) = \frac{1}{\|S\|^2} \sum_{f_i \in S} \sum_{f_j \in S} I(f_i; f_j)$$

- mRMR is **robust** but doesn't measure **interactions** of features in estimating  $c$  (for that we could use ANOVAs).

# 3. $\chi^2$ method

- We adapt the  $\chi^2$  method we saw when testing whether distributions were significantly different:

$$\chi^2 = \sum_{c=1}^C \frac{(O_c - E_c)^2}{E_c} \quad \longrightarrow \quad \chi^2 = \sum_{c=1}^C \sum_{f_i=f}^F \frac{(O_{c,f} - E_{c,f})^2}{E_{c,f}}$$

where  $O_{c,f}$  and  $E_{c,f}$  are the observed and expected number, respectively, of times the class  $c$  occurs together with the (discrete) feature  $f$ .

- The expectation  $E_{c,f}$  assumes  $c$  and  $f$  are **independent**.
- Now, **every feature has a  $p$ -value**. A lower  $p$ -value means  $c$  and  $f$  are *less likely* to be independent.
- Select the  $k$  features with the lowest  $p$ -values.

# Multiple comparisons

- If we're just **ordering** features, this  $\chi^2$  approach is (mostly) fine.
- But what if we get a 'significant'  $p$ -value (e.g.,  $p < 0.05$ )? Can we claim a significant effect of the class on that feature?
- Imagine you're flipping a coin to see if it's fair. You claim that if you get 'heads' in 9/10 flips, it's biased.
- Assuming  $H_0$ , the coin is fair, the probability that a fair coin would come up heads  $\geq 9$  out of 10 times is:

$$(10 + 1) \times 0.5^{10} = 0.0107$$



Number of ways 9 flips are heads    Number of ways all 10 flips are heads

# Multiple comparisons

- But imagine that you're simultaneously testing **173** coins – you're doing **173 (multiple) comparisons**.
- If you want to see if *a specific chosen* coin is fair, you still have only a 1.07% chance that it will give heads  $\geq \frac{9}{10}$  times.
- **But** if you don't preselect a coin, what is the probability that *none* of these fair coins will accidentally appear biased?

$$(1 - 0.0107)^{173} \approx 0.156$$

- If you're testing 1000 coins?

$$(1 - 0.0107)^{1000} \approx 0.0000213$$

# Multiple comparisons

- The more features you evaluate with a statistical test (like  $\chi^2$ ), the more likely you are to accidentally find spurious (incorrect) significance **accidentally**.
- **Bonferroni correction** is an adjustment method:
  - Divide your level of significance required  $\alpha$ , by the number of comparisons.
  - E.g., if  $\alpha = 0.05$ , and you're doing **173** comparisons, each would need  $p < \frac{0.05}{173} \approx 0.00029$  to be considered significant.



# Reading

- Manning & Schütze: 2.2, 5.3-5.5
- Cover & Thomas *Elements of Information Theory*, Chapter 2