

#### **Entropy and decisions**

- HAN THE AR



CSC401/2511 – Natural Language Computing – Spring 2019 Lecture 3, Frank Rudzicz and Chloé Pou-Prom University of Toronto

## **This lecture**

- Information theory and entropy.
- Decisions.
  - Classification.
  - Significance.

Can we quantify the statistical structure in a model of communication? Can we quantify the <u>meaningful</u> 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

Star Trek is better than Star Wars



## 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 *E* when you observe it?



## 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.





#### **Information is additive**

• From kindependent, equally likely events E,

$$I(E^{k}) = \log_{2} \frac{1}{P(E^{k})} = \log_{2} \frac{1}{P(E)^{k}} \qquad I(k \text{ binary decisions}) = \log_{2} \frac{1}{\left(\frac{1}{2}\right)^{k}} = \frac{k \text{ bits}}{\left(\frac{1}{2}\right)^{k}}$$

- For a unigram model, with each of 50K words *w* equally likely,  $I(w) = \log_2 \frac{1}{\frac{1}{\sqrt{50000}}} \approx 15.61 \text{ bits}$ 
  - and for a **sequence** of 1K words in that model,

$$I(w^k) = \log_2 \frac{1}{\left(\frac{1}{50000}\right)^{1000}} \approx ???$$



## Information with unequal events

 An information source S emits symbols without memory from a vocabulary {w<sub>1</sub>, w<sub>2</sub>, ..., w<sub>n</sub>}. Each symbol has its own probability {p<sub>1</sub>, p<sub>2</sub>, ..., p<sub>n</sub>}



Darkside (0.06) Destiny (0.07)

- What is the <u>average</u> amount of information we get in **observing** the **output** of source S ?
  - You still have 6 events that are possible but you're fairly sure it will be 'No'.



#### Entropy

• Entropy: *n*. the average amount of information we get in observing the output of source *S*.

$$H(S) = \sum_{i} p_{i}I(w_{i}) = \sum_{i} p_{i}\log_{2}\frac{1}{p_{i}}$$
ENTROPY

Note that this is *very* similar to how we define the expected value (i.e., 'average') of something:

$$E[X] = \sum_{x \in X} p(x) x$$





CSC401/2511 – Winter 2019

#### **Entropy – examples**



$$H(S) = \sum_{i} p_i \log_2 \frac{1}{p_i}$$
  
= 0.7 log<sub>2</sub>(1/0.7) + 0.1 log<sub>2</sub>(1/0.1) + ...  
= 1.542 bits

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

$$H(S) = \sum_{i} p_{i} \log_{2} \frac{1}{p_{i}} = 6 \left( \frac{1}{6} \log_{2} \frac{1}{1/6} \right)$$
  
= 2.585 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 word, e.g., we'd like a distribution with lower entropy.
  - Low entropy  $\equiv$  less uncertainty



#### **Bounds on entropy**

• Maximum: uniform distribution  $S_1$ . Given M choices,

$$H(S_1) = \sum_{i} p_i \log_2 \frac{1}{p_i} = \sum_{i} \frac{1}{M} \log_2 \frac{1}{1/M} = \log_2 M$$

• Minimum: only one choice,  $H(S_2) = p_i \log_2 \frac{1}{p_i} = 1 \log_2 \frac{1}{p_i} = 0$ 





## **Coding symbols efficiently**

- 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	Huffman Code	
No	000	0	
Yes	001	11	
Destiny	010	101	
Darkside	011	1001	
Maybe	100	10000	
Sure	101	10001	



## **Coding symbols efficiently**

 Another way of looking at this is through the (binary) Huffman tree (*r*-ary trees are often flatter, all else being equal):



Word (sorted)	Linear Code	Huffman Code	
No	000	0	
Yes	001	11	
Destiny	010	101	
Darkside	011	1001	
Maybe	100	10000	
Sure	101	10001	



## **Alternative notions of entropy**

- Entropy is **equivalently**:
  - The average amount of information provided by symbols in a vocabulary,
  - The average amount of uncertainty you have before observing a symbol from a vocabulary,
  - The average amount of 'surprise' you receive when observing a symbol,
  - The number of bits needed to communicate that alphabet
    - Aside: Shannon showed that you cannot have a coding scheme that can communicate the vocabulary more efficiently than H(S)



## **Entropy of several variables**

- Joint entropy
- Conditional entropy
- Mutual information



## **Entropy of several variables**



- Consider the vocabulary of a meteorologist describing
   <u>Temperature and</u> <u>Wetness</u>.
  - <u>T</u>emperature = {hot, mild, cold}
  - <u>W</u>etness = {*dry, wet*}

$$P(W = dry) = 0.6,$$
  
 $P(W = wet) = 0.4$ 
 $H(W) = 0.6 \log_2 \frac{1}{0.6} + 0.4 \log_2 \frac{1}{0.4} = 0.970951$  bits

$$P(T = hot) = 0.3,$$
  
 $P(T = mild) = 0.5,$   
 $P(T = cold) = 0.2$ 

$$H(T) = 0.3 \log_2 \frac{1}{0.3} + 0.5 \log_2 \frac{1}{0.5} + 0.2 \log_2 \frac{1}{0.2} = 1.48548 \text{ bits}$$
  
But W and T are not independent,

Example from Roni Rosenfeld

 $\neq P(VV)P$ 



## 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)

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

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)





## **Entropy given knowledge**

• Consider **conditional** probability, P(T|W)

P(T   W)	T = cold	mild	hot	
W = dry	1/6	2/3	1/6	1.0
wet	1/2	1/4	1/4	1.0

- $H(T|W = dry) = H\left(\left\{\frac{1}{6}, \frac{2}{3}, \frac{1}{6}\right\}\right) = 1.25163$  bits
- $H(T|W = wet) = H\left(\left\{\frac{1}{2}, \frac{1}{4}, \frac{1}{4}\right\}\right) = 1.5$  bits
- Conditional entropy combines these: H(T|W) 0.6 = [p(W = dry)H(T|W = dry)] + [p(W = wet)H(T|W = wet)] = 1.350978 bits



#### **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 = cold	mild	hot
W = dry	0.1/ <mark>0.3</mark>	0.4/ <mark>0.5</mark>	0.1/0.2
wet	0.2/ <mark>0.3</mark>	0.1/ <mark>0.5</mark>	0.1/0.2
	1.0	1.0	1.0

- $H(W|T = cold) = H\left(\left\{\frac{1}{3}, \frac{2}{3}\right\}\right) = 0.918295$  bits
- $H(W|T = mild) = H\left(\left\{\frac{4}{5}, \frac{1}{5}\right\}\right) = 0.721928$  bits
- $H(W|T = hot) = H\left(\left\{\frac{1}{2}, \frac{1}{2}\right\}\right) = 1$  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 hits
- $H(T) H(T|W) \approx 0.1345$  bits

Previously computed

- How much does *T* tell us about *W* on average?
   *H*(*W*) − *H*(*W*|*T*) = 0.970951 − 0.8364528
   ≈ 0.1345 bits
  - Interesting ... is that a coincidence?



## **Mutual information**

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

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

- **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  $\log_2 \frac{p(x,y)}{p(x)p(y)} = \log_2 1 = 0 \ \forall x, y - \text{there is no mutual information}!$



#### **Relations between entropies**





## **Preview – the noisy channel**

Messages can get distorted when passed through a noisy conduit – <u>how much information is lost/retained</u>?



#### **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?





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





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

• 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<sub>i</sub> 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<sub>i</sub> from P, we need log P(i)/O(i) more bits.



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

• Aside:

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



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





## **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.





CSC401/2511 - Winter 2019

#### **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 *p*.
- We estimate p by estimating the probability of a test corpus C using the distribution q:

 $P_q(C)$ 



## **Probability of a corpus?**

 The probability P(C) of a corpus C requires similar assumptions that allowed us to compute the probability P(s<sub>i</sub>) of a sentence s<sub>i</sub>.

	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<sub>i</sub>), 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:  $PP_M(C) = 2^{H(C;M)}$ 



#### 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<sub>0</sub>, which states that the two distributions are <u>the same</u> (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*<sub>0</sub> if it is **too improbable**.
  - How do we determine the value of 'too'?
- Significance level  $\alpha$  ( $0 \le \alpha \le 1$ ) 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% (α = 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:
  - 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  $(\overline{x})$  and variance  $(\sigma)$  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  $H_0$ .

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

An example would be nice...



#### **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 much worse 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 much worse McGill tweets?
- Degrees of freedom (d.f.): n.pl. In this t-test, this is the sum of the number of observations in each group, minus 2 (because there are two groups).
- In our example, we have  $N_{UofT} = 200$  for DCS students, but because we don't sample at McGill,  $N_{McGill} \approx \infty$ , so  $d.f. = \infty$ .
  - (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 much worse 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.
  - 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**.

	lpha (one-tail)	0.05	0.025	0.01	0.005	0.001	0.0005
	1	6.314	12.71	31.82	63.66	318.3	636.6
df	10	1.812	2.228	2.764	3.169	4.144	4.587
u.i.	20	1.725	2.086	2.528	2.845	3.552	3.850
	$\infty$	1.645	1.960	2.326	2.576	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 *α* is ad-hoc.
   There are ways to correct for the selection of *α*.

	lpha (one-tail)	0.05	0.025	0.01	0.005	0.001	0.0005
	1	6.314	12.71	31.82	63.66	318.3	636.6
df	10	1.812	2.228	2.764	3.169	4.144	4.587
u.i.	20	1.725	2.086	2.528	2.845	3.552	3.850
	$\infty$	1.645	1.960	2.326	2.576	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<sub>0</sub> 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<sub>0</sub> is that they are independent, i.e.,

H<sub>0</sub>:  $P(new \ companies) = P(new)P(companies)$ =  $\frac{15828}{14307668} \times \frac{4675}{14307668}$ 

$$\approx 3.615 \times 10^{-7}$$



## 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  $\bar{x} = p = P(new \ companies)$
  - 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

• 
$$\bar{x} = \frac{8}{14307667} \approx 5.591 \times 10^{-7}$$
  
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<sub>0</sub> 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!).



## 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).
- E.g., we measure the accuracy of an ASR system for different settings of empirical parameters M and Q (more on these later in the course...).

Accuracy (%)	M = 2	M = 4	M = 16	$H_0$ : no effect of source variables.			<mark>ibles.</mark>
Q = 2	53.33	66.67	53.33	Source	<b>d</b> . <b>f</b> .	p value	
	26.67	53.33	40.00	0	1	0 179	Accept $H_{0}$
	0.00	40.00	26.67	M	-	0.106	Accept H
Q = 5	93.33	26.67	100.00	1/1	2	0.100	
	66.67	13.33	80.00	interaction	2	0.006	Reject $H_0$ at $\alpha = 0.01$
	40.00	0.00	60.00	A completely fic	tional exa	ample	
		_	_				dia.

#### CSC401/2511 - Winter 2019

# **3.** Pearson's $\chi^2$ test (details aside)

- 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 of and expected number of observations of type c, respectively.



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



- For example, is our die from Lecture 2 fair or not?
- Imagine we throw it 60 times. The expected number of appearances of each side is 10.

С	<b>0</b> <sub>c</sub>	E <sub>c</sub>	$O_c - E_c$	$(\boldsymbol{O}_c - \boldsymbol{E}_c)^2$	$(\boldsymbol{O}_c - \boldsymbol{E}_c)^2 / \boldsymbol{E}_c$
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 ( $\chi^2$ )	13.4	

• With df = 6 - 1 = 5, the critical value is 11.07<**13.4**, so we throw away  $H_0$ : the die is biased.

We'll see χ<sup>2</sup> again soon...



#### Reading

• Manning & Schütze: 2.2, 5.3-5.5



#### **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...