

Natural Language Processing

Yue Zhang
Westlake University



Chapter 5

Using Information Theory

Contents

- 5.1 The maximum entropy principle
 - 5.1.1 Information and entropy
 - 5.1.2 A naïve maximum entropy model
 - 5.1.3 Maximum entropy model and training data
- 5.2 KL-Divergence, Cross-Entropy and Model Perplexity
 - 5.2.1 KL-divergence
 - 5.2.2 Cross entropy
 - 5.2.3 Model Perplexity
- 5.3 Mutual information
 - 5.3.1 Pointwise mutual information
 - 5.3.2 Using PMI

Contents

- 5.1 The maximum entropy principle
 - **5.1.1 Information and entropy**
 - 5.1.2 A naïve maximum entropy model
 - 5.1.3 Maximum entropy model and training data
- 5.2 KL-Divergence, Cross-Entropy and Model Perplexity
 - 5.2.1 KL-divergence
 - 5.2.2 Cross entropy
 - 5.2.3 Model Perplexity
- 5.3 Mutual information
 - 5.3.1 Pointwise mutual information
 - 5.3.2 Using PMI

What is information?

- Resolve uncertainty about random events.

Tossing a coin:



head
"1"



tail
"0"

bit

Draw a card:



diamond
"00"



club
"01"



spade
"10"



heart
"11"

bit

- To learn the outcome of a random event with n equally possible results, $\log_2 n$ bits of information is necessary.

Information of event.

- Original uncertainty ---- remaining uncertainty

Draw a card:



- Spade Ace: $\log_2 52 - \log_2 1 = \log_2 52$
- Spade: $\log_2 52 - \log_2 13 = 2$ bits
- Ace: $\log_2 52 - \log_2 4 = \log_2 13$ bits



m red balls and n green balls in a box

Red ball:

$$\begin{aligned} & \log_2(m + r) - \log_2 n \\ &= \log_2 \frac{m + n}{n} = \log_2 \frac{1}{P(\text{red})} \end{aligned}$$

- The outcomes with higher probabilities contains less information.
- For a certain outcome r_i
 - Probability: $P(r_i)$
 - Information received: $-\log_2 P(r_i)$.

Entropy

- Entropy analyzes information concerning events by considering all possible outcomes or random variables by considering all possible values.



- The entropy of distribution P is:

$$H(P) = - \sum_{i=1}^n P(r_i) \log_2 P(r_i) = E \left(\log_2 \frac{1}{P(r_i)} \right)$$

- where E denotes a probability-weighted average, or the **mathematical expectation**

Entropy and distribution characteristics

- Encoding six numbers

- Uniform distribution

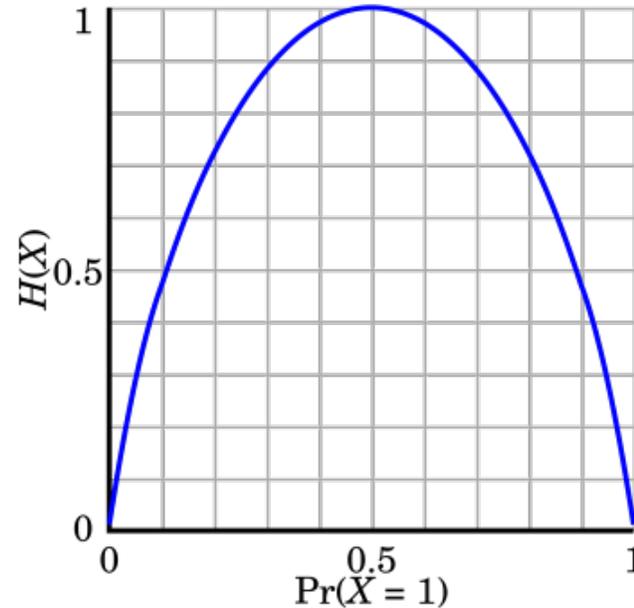
3 bits 001 010 011 100 101 110

- 90% case with number 1, and 2% with 2, 3, 4, 5, 6

0 1010 1011 1110 1101 1110

$0.9 \times 1 + 0.1 \times 4 = 1.3$ bits

Entropy



binary entropy for
tossing a coin

the expected surprisal

$X=1$ represents a result
of heads

Entropy $H(X)$ of a coin flip corresponds to probability $\Pr(X)$

- Events with uniform output distributions have the largest entropy.
- The more uneven the distribution is, the smaller the entropy is.

Contents

- 5.1 The maximum entropy principle
 - 5.1.1 Information and entropy
 - **5.1.2 A naïve maximum entropy model**
 - 5.1.3 Maximum entropy model and training data
- 5.2 KL-Divergence, Cross-Entropy and Model Perplexity
 - 5.2.1 KL-divergence
 - 5.2.2 Cross entropy
 - 5.2.3 Model Perplexity
- 5.3 Mutual information
 - 5.3.1 Pointwise mutual information
 - 5.3.2 Using PMI

Occam's razor

- A principle attributed to the fourteenth-century English Franciscan friar William of Ockham, which states that *“entities should not be multiplied beyond necessity”*
- **Occam's razor** shares underlying similarities with the **maximum entropy principle**.



A naive maximum entropy model

A probabilistic model for a random event e with possible outcomes r_1, r_2, \dots, r_M :

$$\hat{P} = \operatorname{argmax} H(P) = \operatorname{argmax} - \sum_{i=1}^M P(r_i) \log_2 P(r_i)$$

using $P(r_i)$ directly as parameters.

The training objective is to find:

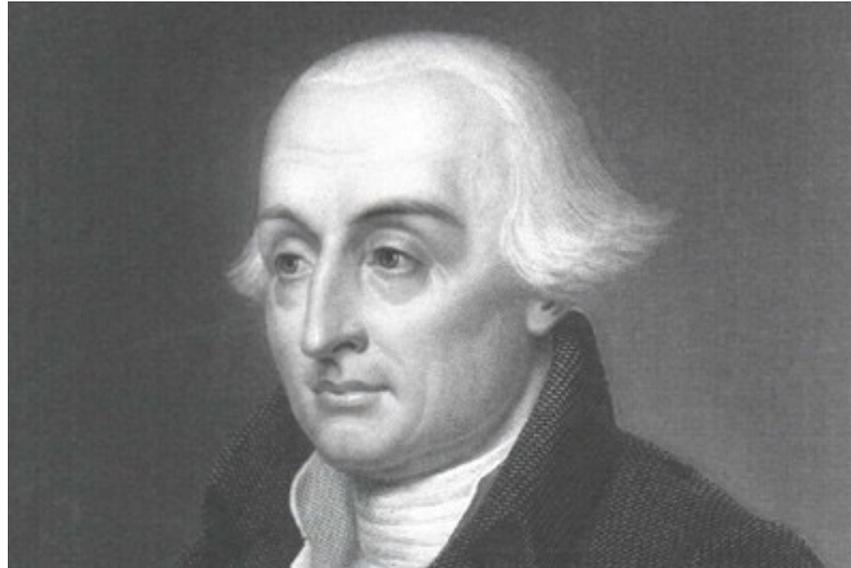
$$\hat{P}(e) = \operatorname{arg max} H(e) = \operatorname{arg min} \sum_{i=1}^M P(r_i) \log_2 P(r_i)$$

under the constraint that

$$\sum_{i=1}^M P(r_i) = 1$$

A naive maximum entropy model

Taking each $P(r_i)$ as a separate variable, we use **Lagrange multipliers** to do the optimization.



In mathematic optimization, the *method of Lagrange multipliers* is a strategy for finding the local maxima and minima of a function subject to equality constraints.

A naive maximum entropy model

- The Lagrangian equation is

$$\Lambda(P(r_1), P(r_2), \dots, P(r_M), \lambda) = \sum_{i=1}^M P(r_i) \log_2 P(r_i) + \lambda(\sum_{i=1}^M P(r_i) - 1),$$

where λ is a Lagrangian multiplier.

- A necessary condition for optimality in the constrained problem is that

$$\frac{\partial \Lambda}{\partial P(r_i)} = 0 \text{ for } i \in [1 \dots m] \quad \Rightarrow \quad 1 - \log_2 P(r_i) + \lambda = 0$$

which suggests that $P(r_1) = P(r_2) = \dots = P(r_M)$

- The conclusion conforms the fact that the uniform distribution contains the most uncertainty.

Contents

- 5.1 The maximum entropy principle
 - 5.1.1 Information and entropy
 - 5.1.2 A naïve maximum entropy model
 - **5.1.3 Maximum entropy model and training data**
- 5.2 KL-Divergence, Cross-Entropy and Model Perplexity
 - 5.2.1 KL-divergence
 - 5.2.2 Cross entropy
 - 5.2.3 Model Perplexity
- 5.3 Mutual information
 - 5.3.1 Pointwise mutual information
 - 5.3.2 Using PMI

Conditional entropy

For a conditional probability distribution $P(y|x)$, given that the random event X follows a probability distribution $P(x)$, the conditional entropy value:

$$\begin{aligned} H(Y|X) &= - \sum_x \sum_y P(x)P(y|x) \log_2 P(y|x) \\ &= - \sum_x \sum_y P(x, y) \log_2 P(y|x) \end{aligned}$$

Conditional entropy

For a conditional probability distribution $P(y|x)$, given that the random event X follows a probability distribution $P(x)$, the conditional entropy value:

$$\begin{aligned} H(Y|X) &= - \sum_x \sum_y P(x) P(y|x) \log_2 P(y|x) \\ &= - \sum_x \sum_y P(x, y) \log_2 P(y|x) \end{aligned}$$

Intuition

- Given x , $H(Y|x) = - \sum_y P(y|x) \log_2 P(y|x)$
- Expectation $H(Y|X) = \sum_x P(x) H(Y|x)$

Maximum entropy model and training data

- We are to derive a maximum entropy model for feature-based discriminative classification.
 - Training data : $D = \{(x_i, y_i)\}_{i=1}^N$
 - Feature instances for $(x, y) : f_1, f_2, \dots, f_m$
 - Feature instance : $f_i(x, y)$ (count)

Notations

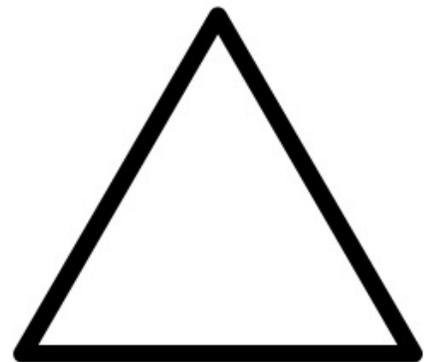
- The model to build : $P(y | x)$
- Prior distribution of x : $P(x)$
- Empirical distribution : $\tilde{P}(x) = \frac{\#x}{\sum_{\{x' \in D\}} \#x'} = \frac{\#x}{|D|}$
- Model expectation of f_i : $E(f_i)$
- Empirical count of f_i : $\tilde{E}(f_i)$

Modelling the problem

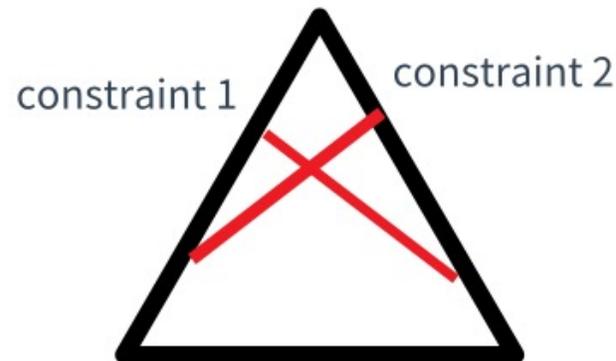
- Objective function : $H(P)$
- Goal : among all distributions that satisfy the constraints, choose the one \hat{P} that maximizes $H(P)$

$$\hat{P} = \operatorname{argmax} H(P)$$

- Constraints: feature counts.



The parameter space
of a probabilistic model



models satisfy
the constraints

Objective

- The conditional entropy to maximize is
$$H(Y|X) = - \sum_x \sum_y P(x)P(y|x) \log_2 P(y|x)$$
- We use $\tilde{P}(x) = \frac{\#x}{\sum_{x' \in D} \#x'} = \frac{\#x}{|D|}$ to represent $P(x)$, which is typically $\frac{1}{|D|}$.
- resulting in $H(Y|X) = - \sum_x \sum_y \tilde{P}(x) P(y|x) \log_2 P(y|x)$, which we maximize

Constraints

- Model's feature expectation = observed feature expectation $E(f_i) = \tilde{E}(f_i)$

- The *model* feature

$$E(f_i) = \sum_{j=1}^{|D|} \tilde{P}(x_j) \sum_y P(y|x_j) f_i(x_j, y)$$

typically $\frac{1}{|D|} \sum_{j=1}^{|D|} \sum_y P(y|x_j) f_i(x_j, y)$.

- The *empirical* feature

$$\tilde{E}(f_i) = \sum_{j=1}^{|D|} \tilde{P}(x_j, y_j) f_i(x_j, y_j) \quad \left((x_j, y_j) \in D \right)$$

typically $\frac{1}{|D|} \sum_{j=1}^{|D|} f_i(x_j, y_j)$.

- One additional constraint, as before $\sum_y P(y|x) = 1$.

Using Lagrangian multipliers

- A model $\hat{P}(y|x)$ that satisfies

$$\hat{P}(y|x) = \arg \min \sum_x \sum_y \tilde{P}(x) P(y|x) \log_2 P(y|x)$$

s.t. for all $i, E(f_i) = \tilde{E}(f_i); \sum_y P(y|x) = 1$

- The Lagrangian equation is

$$\Lambda(P, \vec{\lambda}) = -H(Y|X) + \sum_{i=1}^m \lambda_i (E(f_i) - \tilde{E}(f_i)) + \sum_x \lambda_{m+1}^x \left(\sum_y P(y|x) - 1 \right)$$

- A necessary condition to the constrained minimum value of $-H(Y | X)$ is $\frac{\partial \Lambda}{\partial P} = 0$

- Solving these equations, we have

$$P(y|x) = \frac{\exp(\sum_i \lambda_i f_i(x,y))}{\sum_{y'} \exp(\sum_i \lambda_i f_i(x,y'))}$$

- This is a log-linear form of $P(y | x)$ using the maximum entropy principle, which is the same as a log linear model
- We further find $\vec{\lambda}$ via $\vec{\lambda} = \arg \min_{\vec{\lambda}} \Lambda(P, \lambda) = \arg \min_{\vec{\lambda}} - \sum_j P(y_i | x_i)$, which is the same as MLE.

Contents

- 5.1 The maximum entropy principle
 - 5.1.1 Information and entropy
 - 5.1.2 A naïve maximum entropy model
 - 5.1.3 Maximum entropy model and training data
- 5.2 KL-Divergence, Cross-Entropy and Model Perplexity
 - **5.2.1 KL-divergence**
 - 5.2.2 Cross entropy
 - 5.2.3 Model Perplexity
- 5.3 Mutual information
 - 5.3.1 Pointwise mutual information
 - 5.3.2 Using PMI

KL-divergence

- **Revisit risk.** Parameter $\vec{\theta}$, data $D = \{d_i\}_{i=1}^N$

$$\overline{risk}(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N \text{loss}(\vec{\theta} \cdot \vec{\phi}(d_i))$$

KL-divergence

- **Revisit risk.** Parameter $\vec{\theta}$, data $D = \{d_i\}_{i=1}^N$

$$\widetilde{risk}(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N \text{loss}(\vec{\theta} \cdot \vec{\phi}(d_i))$$

- **For a probabilistic model**

$$\widetilde{risk}(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N \text{diff}(\tilde{P}(d_i), Q(d_i))$$

where $\tilde{P}(d_i)$ is data frequency , $Q(d_i)$ is model probability.

KL-divergence

- **Revisit risk.** Parameter $\vec{\theta}$, data $D = \{d_i\}_{i=1}^N$

$$\widetilde{risk}(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N \text{loss}(\vec{\theta} \cdot \vec{\phi}(d_i))$$

- **For a probabilistic model**

$$\widetilde{risk}(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N \text{diff}(\tilde{P}(d_i), Q(d_i))$$

where $\tilde{P}(d_i)$ is data frequency, $Q(d_i)$ is model probability.

- ***diff* can be defined**

$$\widetilde{risk}(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N (\log_2 \tilde{P}(d_i) - \log_2 Q(d_i))$$

$$= \sum_{i=1}^N \tilde{P}(d_i) (\log_2 \tilde{P}(d_i) - \log_2 Q(d_i)) = \sum_{i=1}^N \tilde{P}(d_i) \log_2 \frac{\tilde{P}(d_i)}{Q(d_i)}$$

KL-divergence

- **Revisit risk.** Parameter $\vec{\theta}$, data $D = \{d_i\}_{i=1}^N$

$$\widetilde{risk}(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N \text{loss}(\vec{\theta} \cdot \vec{\phi}(d_i))$$

- **For a probabilistic model**

$$\widetilde{risk}(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N \text{diff}(\tilde{P}(d_i), Q(d_i))$$

where $\tilde{P}(d_i)$ is data frequency, $Q(d_i)$ is model probability.

- ***diff* can be defined**

$$\widetilde{risk}(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N (\log_2 \tilde{P}(d_i) - \log_2 Q(d_i))$$

$$= \sum_{i=1}^N \tilde{P}(d_i) (\log_2 \tilde{P}(d_i) - \log_2 Q(d_i)) = \sum_{i=1}^N \tilde{P}(d_i) \log_2 \frac{\tilde{P}(d_i)}{Q(d_i)}$$

- **Kullback-Leibler (KL) divergence** measures how different two distributions of the same random variable are.

$$KL(P, Q) = \sum_{i=1}^M P(\tau_i) \log_2 \frac{P(\tau_i)}{Q(\tau_i)} = E_{e \sim P(e)} \log_2 \frac{P(e)}{Q(e)}$$

- Not symmetric --- how different is Q according to P .
- Can measure a probabilistic model against a data distribution.
 $KL(P, Q) \geq 0$, $KL(P, Q) = 0$, where $P = Q$.

- **Loss function:**
$$KL(P, Q) = \sum_{i=1}^N \tilde{P}(d_i) \left(\log_2 \tilde{P}(d_i) - \log_2 Q(d_i) \right)$$
$$= \sum_{i=1}^N \tilde{P}(d_i) \log_2 \tilde{P}(d_i) - \sum_{i=1}^N \tilde{P}(d_i) \log_2 Q(d_i)$$

- As the first term is constant, the loss effectively maximizes

$$\sum_{i=1}^N \tilde{P}(d_i) \log_2 Q(d_i) = \frac{1}{N} \sum_{i=1}^N \log_2 Q(d_i)$$

which is the exactly log-likelihood of the dataset D , namely **MLE**.

Contents

- 5.1 The maximum entropy principle
 - 5.1.1 Information and entropy
 - 5.1.2 A naïve maximum entropy model
 - 5.1.3 Maximum entropy model and training data
- 5.2 KL-Divergence, Cross-Entropy and Model Perplexity
 - 5.2.1 KL-divergence
 - **5.2.2 Cross entropy**
 - 5.2.3 Model Perplexity
- 5.3 Mutual information
 - 5.3.1 Pointwise mutual information
 - 5.3.2 Using PMI

Cross entropy

- The second term of KL-divergence is referred to **cross-entropy**:

$$H(P, Q) = - \sum_{i=1}^M P(r_i) \log_2 Q(r_i) = E_{e \sim P} \log_2 \frac{1}{Q(e)}$$

It also measures the similarity between two distributions of the **same** random variable.

Cross entropy

- The second term of KL-divergence is referred to **cross-entropy**:

$$H(P, Q) = - \sum_{i=1}^M P(r_i) \log_2 Q(r_i) = E_{e \sim P} \log_2 \frac{1}{Q(e)}$$

It also measures the similarity between two distributions of the **same** random variable.

- Intuitively, it means the number of bits to encode a variable e distributed in Q using the encoding scheme defined by P .

Thus $H(P, Q) = H(P)$ if $Q = P$ and is larger when Q differs more from P .

Cross entropy

- The second term of KL-divergence is referred to **cross-entropy**:

$$H(P, Q) = - \sum_{i=1}^M P(r_i) \log_2 Q(r_i) = E_{e \sim P} \log_2 \frac{1}{Q(e)}$$

It also measures the similarity between two distributions of the **same** random variable.

- Intuitively, it means the number of bits to encode a variable e distributed in Q using the encoding scheme defined by P .

Thus $H(P, Q) = H(P)$ if $Q = P$ and is larger when Q differs more from P .

- KL-divergence is non-negative because:

$$KL(P, Q) = \sum_{i=1}^N \tilde{P}(d_i) \log_2 \tilde{P}(d_i) - \sum_{i=1}^N \tilde{P}(d_i) \log_2 Q(d_i) = H(P, Q) - H(P)$$

As a result, KL-divergence is also called **relative entropy**.

Cross entropy loss

- Cross-entropy: $H(P, Q) = - \sum_{i=1}^M P(r_i) \log_2 Q(r_i)$
- Cross-entropy loss:

$$H(\tilde{P}, Q) = - \sum_{i=1}^N \tilde{P}(d_i) \log_2 Q(d_i) = - \frac{1}{N} \sum_{i=1}^N \log_2 Q(d_i)$$

where $\tilde{P}(d_i)$ is data frequency, $Q(d_i)$ is model probability.

- The same as negative log-likelihood loss.

Contents

- 5.1 The maximum entropy principle
 - 5.1.1 Information and entropy
 - 5.1.2 A naïve maximum entropy model
 - 5.1.3 Maximum entropy model and training data
- 5.2 KL-Divergence, Cross-Entropy and Model Perplexity
 - 5.2.1 KL-divergence
 - 5.2.2 Cross entropy
 - **5.2.3 Model Perplexity**
- 5.3 Mutual information
 - 5.3.1 Pointwise mutual information
 - 5.3.2 Using PMI

Perplexity

- Formally $\Upsilon(P) = 2^{H(P)} = 2^{-\sum_i P(z_i) \log_2 P(z_i)}$
- Intuitively, **perplexity** represents the *expected* number of bits necessary for encoding each outcome.

- Cross-entropy can also be used as the power term for calculating perplexity.
- This can be useful for model evaluation.

$$\Upsilon(Q, D) = 2^{H(\tilde{P}(d), Q)} = 2^{-\sum_{i=1}^N \tilde{P}(d_i) \log_2 Q(d_i)} = 2^{-\frac{1}{N} \sum_{i=1}^N \log_2 Q(d_i)}$$

(Model perplexity)

where $\tilde{P}(d_i)$ is data frequency, $Q(d_i)$ is model probability.

Evaluating language models

- For classification, accuracy is the metric.
- For language modeling, no single correct answer.
 - For sentence level, $2^{-\frac{1}{N} \sum_{i=1}^N \log_2 Q(s_i)}$, typically 2^{190} .
 - A commonly used evaluation metric for language models is **per-word perplexity**:

$$2^{-\frac{1}{|D|} \sum_{i=1}^{|D|} \log Q(w_i)}$$

typically 10-250.

Contents

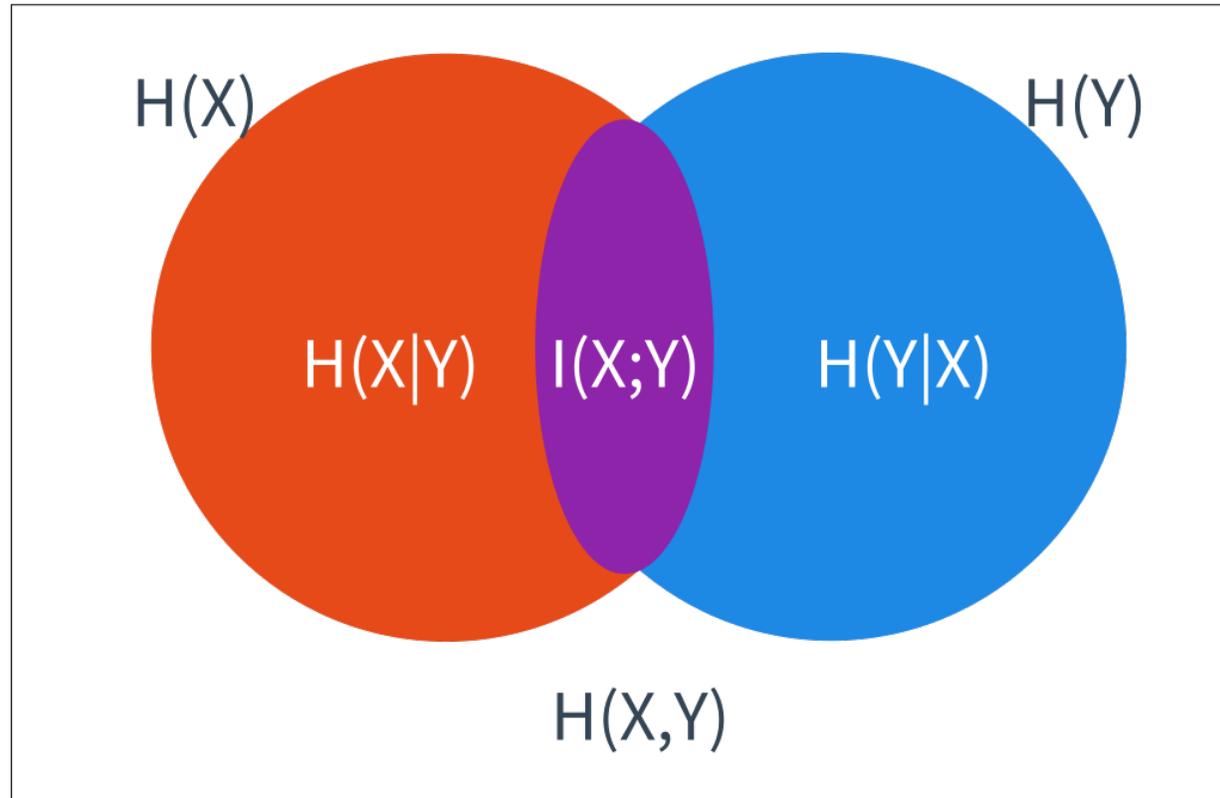
- 5.1 The maximum entropy principle
 - 5.1.1 Information and entropy
 - 5.1.2 A naïve maximum entropy model
 - 5.1.3 Maximum entropy model and training data
- 5.2 KL-Divergence, Cross-Entropy and Model Perplexity
 - 5.2.1 KL-divergence
 - 5.2.2 Cross entropy
 - 5.2.3 Model Perplexity
- **5.3 Mutual information**
 - 5.3.1 Pointwise mutual information
 - 5.3.2 Using PMI

- Entropy – the *expected* number of bits to encode a random variable.
The number of bits in average when encoding many outcome values.
Optimal encoding scheme.
- Cross-Entropy – the *expected* number of bits to encode a random variable using a different encoding scheme.
Two distributions concerning the same random variable.
 - Evaluate model distribution against data distribution
 - Calibrate model distribution against data.
- Mutual Information – the *expected* number of bits you can save for encoding a random variable if a second random variable is known.
About two random variables.

Mutual information

- Measures the correlation between two different random variables X and Y .
- The difference between $H(Y)$ and $H(Y|X)$ is called the **mutual information** between X and Y , denoted as $I(X, Y)$
- $$I(X, Y) = H(Y) - H(Y|X)$$
$$= \sum_{x,y} P(x, y) \log_2 P(y|x) - \sum_y P(y) \cdot \log_2 P(y)$$
$$= \sum_{x,y} P(x, y) \log_2 \frac{P(x,y)}{P(x)P(y)}$$
- It measures the number of bits we can save for encoding Y , if X is known.

Conditional entropy



Contents

- 5.1 The maximum entropy principle
 - 5.1.1 Information and entropy
 - 5.1.2 A naïve maximum entropy model
 - 5.1.3 Maximum entropy model and training data
- 5.2 KL-Divergence, Cross-Entropy and Model Perplexity
 - 5.2.1 KL-divergence
 - 5.2.2 Cross entropy
 - 5.2.3 Model Perplexity
- 5.3 Mutual information
 - **5.3.1 Pointwise mutual information**
 - 5.3.2 Using PMI

Pointwise mutual information

- Given two random events X and Y , their mutual information can be viewed as the expectation of $\log_2 \frac{P(x,y)}{P(x)P(y)}$ over all x, y :

$$I(x, y) = \sum_{x,y} P(x, y) \log_2 \frac{P(x, y)}{P(x)P(y)} = E_{x,y} \left(\log_2 \frac{P(x, y)}{P(x)P(y)} \right)$$

- For each outcome pair (x,y) , $\log_2 \frac{P(x,y)}{P(x)P(y)}$ is called Pointwise Mutual information (PMI) between x and y .

Pointwise mutual information

$$\log_2 \frac{P(x, y)}{P(x)P(y)}$$

- PMI represents the statistical correlation between two values of a random variable, or two outcomes of a random event.
- PMI --- Mutual information
information --- entropy
- PMI can be negative!

Contents

- 5.1 The maximum entropy principle
 - 5.1.1 Information and entropy
 - 5.1.2 A naïve maximum entropy model
 - 5.1.3 Maximum entropy model and training data
- 5.2 KL-Divergence, Cross-Entropy and Model Perplexity
 - 5.2.1 KL-divergence
 - 5.2.2 Cross entropy
 - 5.2.3 Model Perplexity
- 5.3 Mutual information
 - 5.3.1 Pointwise mutual information
 - **5.3.2 Using PMI**

- **The correlation between variables.**
 - Words and sentiment signals.
 - Neighboring words.
 - Features and class labels.

- **Sentiment lexicon** contains information about the polarity and strength of sentiment words.
- $LEX(w)$ represents the sentiment polarity, and the absolute value represents the strength.

$$SENTI(d) = \frac{\sum_i Lex(w_i)}{|\{w_i | Lex(w_i) \neq 0\}|}$$

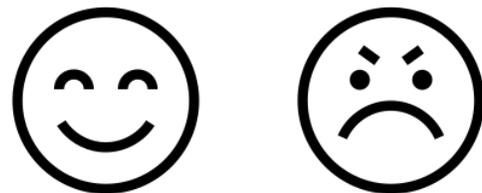
Learning sentiment lexicons

- The PMI between a word w and a *seed* word

$$PMI(w, seed) = \log_2 \frac{P(w, seed)}{P(w)p(seed)}$$

$$LEX(w) = PMI(w, good) - PMI(w, bad)$$

Emoticons in social media



Collocation extraction

- *Collocation* refers to words that are conventionally used together for certain meaning.



Mr President
Mr Executive



High temperature
Big temperature

- Given two words w_1 and w_2 and a corpus D , their association can be calculated using

$$PMI(w_1, w_2) = \log_2 \frac{P(w_1 w_2)}{P(w_1)P(w_2)}$$

Using PMI to select features

- Feature selection: reduce the size of the feature vector

important features

unimportant features

goal , statement , president

a , in , does

- PMI between feature and class is a commonly used metric for feature selection, the higher PMI value, the more likely w is a strong indicator of c .

- Representing a word in vector space
 - Useful for measuring semantic correlations
 - Thus far we have only learned a “one-hot” representation

PMI and vector representations of words

- *Distributional semantics* : the company of a word (*k-word windows*) tells us much information about its attributes.

Sentence	k	Context
s_1	2	{between, the, and, the }
	5	{the, water, halfway, between, the, and, the, island, ., <s>}
	7	{out, of, the, water, halfway, between, the, and, the, island, ., <s>, <s>, <s>}
s_2	2	{with, the, statement, are}
	5	{are, not, included, with, the, statement, are , called, outstanding, checks}
	7	{written but, are, not, included, with, the, statement, are , called, outstanding, checks, . , <s>}

- K-word windows for the word "bank" in
 s_1 --- *There happened to be a rock sticking out of the water halfway between the bank and the island*
 s_2 --- *The checks that have been written but are not included with the bank statement are called outstanding checks*

Word Representation

Word	Representation	Feature vector
<i>cat</i>	One-hot	$\langle f_1 = 0, \dots, f_{121} = 1, \dots, f_{500} = 0, \dots, f_{10000} = 0 \rangle$
	Context	$\langle f_1 = 1280, f_2 = 0, \dots, f_{35} = 332, \dots, f_{10000} = 0 \rangle$
<i>dog</i>	One-hot	$\langle f_1 = 0, \dots, f_{121} = 0, \dots, f_{500} = 1, \dots, f_{10000} = 0 \rangle$
	Context	$\langle f_1 = 1190, f_2 = 19, \dots, f_{35} = 271, \dots, f_{10000} = 0 \rangle$

cat --- 121, dog --- 500, considering vector dot product.

- The vector representation of a word w_i is :

$$\overline{Vec(w_i)} = \langle PMI(w, w_1), PMI(w, w_2), \dots, PMI(w, w_{|V|}) \rangle$$

$$\begin{aligned} P(u, v) &= \frac{\#(u \text{ and } v \text{ in each other's context window})}{\#(\text{any two words in each other's context window})} \\ &= \frac{\#(u \text{ and } v \text{ in each other's context window})}{2k|D|^2}. \end{aligned}$$

Word Representation

Word	Representation	Feature vector
<i>cat</i>	One-hot	$\langle f_1 = 0, \dots, f_{121} = 1, \dots, f_{500} = 0, \dots, f_{10000} = 0 \rangle$
	Context	$\langle f_1 = 1280, f_2 = 0, \dots, f_{35} = 332, \dots, f_{10000} = 0 \rangle$
	PPMI	$\langle f_1 = 0.3, f_2 = 0, \dots, f_{35} = 2.32, \dots, f_{10000} = 0 \rangle$
<i>dog</i>	One-hot	$\langle f_1 = 0, \dots, f_{121} = 0, \dots, f_{500} = 1, \dots, f_{10000} = 0 \rangle$
	Context	$\langle f_1 = 1190, f_2 = 19, \dots, f_{35} = 271, \dots, f_{10000} = 0 \rangle$
	PPMI	$\langle f_1 = 0.44, \dots, f_{12} = 0.05, \dots, f_{35} = 5.56, \dots, f_{10000} = 0 \rangle$

- PMI and TF-IDF.

- $PPMI(u, v) = \max(PMI(u, v), 0)$
- We use positive PMI (PPMI) to reduce noise and the non-informative.
- Using PPMI, non differentiating words will have a small contribution to the distributional word representation.

Summary

- Entropy and information.
- The maximum entropy principle for defining probabilistic models, and its application in deriving log-linear model forms
- Model perplexity, cross-entropy and KL-divergence for measuring the consistence between model distributions and data distributions
- Mutual information and pointwise mutual information (PMI) for natural language tasks
- Word representations and pointwise mutual information.