

Language ambiguity (has multiple precise meanings): - word 'bring me the file'(resolve w/ POS) - syntactic 'I shot an elephant in my pjs' - semantic 'the rabbit is ready for lunch' - referential 'Pavarotti is a big opera star' - non-literal 'it's raining cats and dogs'		Lexicon – morph(eme)ological analysis (stem and affix e.g. 'cat'+ 's') Deep Learning learns/abstract functions instead of rules based on maintenance of intuitive linguistic rules.	Word segmentation (tokenization) Word normalization (case/acronyms/spelling) Lemmatization 'sing, sung, sang' → 'sing' Stemming (common root, above 's') Part-Of-Speech (tag words with noun, verb...)	Context-Free Grammar: Derive sentence structure through a parse tree S → NP VP, NP → Det N, VP → V NP, VP → V, VP → V PP, PP → P NP <u>Discourse</u> : meaning of a text (relationship between sentences) <u>Pragmatics</u> : intentions/commands <u>Corpus</u> : a collection of documents <u>Document</u> : one item of corpus (sequence) <u>Token</u> : atomic word unit <u>Vocabulary</u> : unique tokens across corpus. One-Hot Encoding : sparse (wasted space), orthogonal vectors (every word is equidistant), cannot represent out of vocab well
Sigmoid (binary class.) $\frac{1}{1+e^{-x}}$, ReLU: max(0, x), Tanh: $\frac{e^x - e^{-x}}{e^x + e^{-x}}$, Softmax (k-class): $\frac{e^{x_k}}{\sum_{k=1}^K e^{x_k}}$ MSE (regression) $\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$, Binary cross-entropy: $-\frac{1}{N} \sum_{i=1}^N (y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}))$ Categorical cross entropy (k-class): $-\frac{1}{N} \sum_{i=1}^N \sum_{c=1}^C y_c^{(i)} \log(\hat{y}_c^{(i)})$		Euclidean Distance: $\sqrt{\sum_{i=1}^n (q_d - d_i)^2}$ Cosine Similarity: $\cos(\theta) = \frac{p_1 \cdot p_2}{ p_1 \times p_2 }$ Analogy Recovery : offset of the vectors reflect their relationship. $a - b \approx c - d \iff d \approx c - a + b$	Window : window consists of target and context (surrounding), Window size = radius Continuous Bag Of Words : context → target, Skip-gram : target → context (give as one-hot, get word representation, map embedding to target words using weight matrix, apply softmax). Train with list of pairs (target, context) by sliding window over input. Loss : $p(w_{t+j} w_t) = \frac{\exp(u_{w_{t+j}} \cdot h_{w_t})}{\sum_{w'=1}^V \exp(u_{w'} \cdot h_{w_t})}$, the aim: $\max \prod_t \prod_j p(w_{t+j} w_t) \rightarrow \min_{\theta} - \sum_t \sum_j \log p(w_{t+j} w_t; \theta) \rightarrow \frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j} w_t; \theta)$ over all elems in corpus. However, the bottom term in the $p(w_{t+j} w_t)$ is inefficient to compute across the entire corpus vocabulary. Therefore, train a Negative Sampling model to predict whether a word appears in the context of another: $\log p(D = 1 w_t, w_{t+1}) + k \mathbb{E}_{\tilde{c} \sim P_{\text{context}}} [\log p(D = 0 w_t, \tilde{c})]$ where $p(D = 1 w_t, w_{t+1})$ is a binary logistic regression probability of seeing the word w_t in the context w_{t+1} . Approximate the expectation by drawing random words from vocabulary, and on left choose positive pairs. Thus the equation is replaced: $p(D = 1 w_t, w_{t+1}) = \frac{1}{1 + \exp(-u_{w_{t+1}} \cdot h_{w_t})}$. We can sample k (5-10 words) with frequency or random sampling.	
Classification : $\hat{y} = \arg \max_y P(y x)$ predict which y is most likely given input x. In the MultiNLI corpus we are given pairs of sentences (premise , hypothesis) with classification problem (Entailment): If hypothesis is implied by premise, Contradiction : If hypothesis contradicts the premise, Neutral : otherwise).		$P(y x) = \frac{\overbrace{P(x y)}^{\text{Likelihood}} \overbrace{P(y)}^{\text{Prior}}}{\underbrace{P(x)}_{\text{Evidence}}}$ Naive Bayes Classifier : $\hat{y} = \arg \max_y \underbrace{P(x_1 y) \dots P(x_l y)}_{P(x_1 y) \dots P(x_l y)} P(y) = \arg \max_y P(y) \prod_{i=1}^l P(x_i y)$		
<i>Accuracy</i> = $\frac{TP+TN}{TP+FP+TN+FN}$, $f1 = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}} = \frac{TP}{TP+0.5(FP+FN)}$, Macro average: averaging of each class F1 scores: increases the emphasis on less frequent classes. Micro average: TPs, TNs, FNs. FPs are summed across each class e.g. $\frac{\sum_i^C TP_i}{\sum_i^C TP_i + \frac{1}{2}(\sum_i^C FP_i + \sum_i^C FN_i)} = \text{Accuracy}$		In a Bag of Words count the number of times a token appears in the vocabulary per class. In One smoothed Naive Bayes : $P(x_i y) = \frac{\text{count}(x_i, y) + 1}{\sum_{x \in V} (\text{count}(x, y) + 1)} = \frac{\text{count}(x_i, y) + 1}{(\sum_{x \in V} \text{count}(x, y)) + V }$ Binary Naive Bayes : only consider if a feature is present, rather than considering every time it occurs. Controlling for negation : pre-pend 'NOT_'.		
Language Models : Assign probabilities to sequence of words, like predicting the next word in a sentence. Uni-directional : use information from left to generate predictions about words on the right. Bi-directional : use information from both sides to fill the target. N-gram : need because language is flexible, and a natural extension of a sentence may no appear in corpus. $P(w_n w_1^{n-1}) \approx P(w_n w_{n-N+1}^{n-1}) = \frac{C(w_{n-N+1}^{n-1} w_n)}{C(w_{n-N+1}^{n-1})}$. To make a prediction $P(w_1, \dots, w_n) = \prod_{k=1}^n P(w_k w_1^{k-1})$ we switch into log space $\log P[\cdot] = \sum_{k=1}^n \log(P(w_k w_1^{k-1}))$ however the longer the sentence the lower its likelihood. Switch to Perplexity : where n is the number of words: $PPL(w) = \sqrt[n]{\prod_{i=1}^n P(w_i w_1^{i-1})}$ the higher the conditional probability of the word sequence, the lower the perplexity. Thus, minimizing perplexity is equivalent to maximizing the test set probability according to the language model. It is a measure of surprise in an LM when seeing new text. For a single word, the score is 1. If the goal of the language model is to support with another task, the best choice of language model is the one that improves downstream task performance the most (extrinsic evaluation). Perplexity is less useful in this case (intrinsic evaluation).		Discriminative algorithms directly learn P(Y X) without considering likelihood. Generative : consider likelihood Logistic Regression : apply sigmoid/softmax with $s = w \cdot x + b$ with loss $H(P, Q) = -\sum_i P(y_i) \log Q(y_i)$.		
Cross Entropy Loss : The cross-entropy is useful when we don't know the actual probability distribution p that generated some data. $\hat{\theta}_{ML} = \arg \max_{\theta} p_{model}(Y X; \theta)$ $= \arg \max_{\theta} \frac{1}{N} \prod_{i=1}^N p_{model}(y^{(i)} x^{(i)}; \theta)$ $= \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N \log p_{model}(y^{(i)} x^{(i)}; \theta)$ $= \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N -\log p_{model}(y^{(i)} x^{(i)}; \theta)$ $= \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N -\log p_{model}(y^{(i)} x^{(i)}; \theta) - \log p(x^{(i)})$ $= \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N -\log p_{model}(y^{(i)} x^{(i)}; \theta) - \log p_{model}(x^{(i)} \theta)$ $= \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N -\log (p_{model}(y^{(i)} x^{(i)}; \theta)p_{model}(x^{(i)} \theta))$ $= \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N -\log p_{model}(y^{(i)}, x^{(i)} \theta)$ $\approx \arg \min_{\theta} (\mathbb{E}_{p_{data}(x,y)} [-\log p_{model}(y, x \theta)])$ $= \arg \min_{\theta} (\mathbb{E}_{p_{data}(x,y)} [-\log p_{model}(\theta)(y, x)])$		RNN : $h_{t+1} = f(h_t, x_t) = \tanh(W h_t + U x_t)$, $W \in \mathbb{R}^{H \times H}$, $U \in \mathbb{R}^{H \times E}$ The model is less able to learn from earlier inputs, better for long ranged CNN : CNNs can perform well if the task involves key phrase recognition		