

NLP Review

Lara J. Martin (she/they)

TA: Omkar Kulkarni (he)

<https://laramartin.net/NLP-class/>

Instructions for Today's Review

1. Get in groups of about 4 people
2. Go through each person's questions and teach each other
3. Make a note of any questions that couldn't be answered, things you still need clarity on, or new questions that came up
4. We'll discuss those questions as a class in the second half of class time

Classification Evaluation: Accuracy, Precision, and Recall

Min: 0 😞

Max: 1 😊

Accuracy: % of items correct

$$\frac{TP + TN}{TP + FP + FN + TN}$$

Precision: % of selected items that are correct

$$\frac{TP}{TP + FP}$$

Recall: % of correct items that are selected

$$\frac{TP}{TP + FN}$$

“**Recall** measures the percentage of items actually present in the input that were correctly identified by the system.”
SLP, ch. 4

	Actually Target	Actually Not Target
Selected/Guessed	True Positive (TP)	False Positive (FP)
Not select/not guessed	False Negative (FN)	True Negative (TN)

F1

$$F_1 = \frac{2 * P * R}{P + R} = \frac{2 * TP}{2 * TP + FP + FN}$$

P/R/F in a Multi-class Setting: Micro- vs. Macro-Averaging

Macroaveraging: Compute performance for each class, then average.

$$\text{macroprecision} = \frac{1}{C} \sum_c \frac{TP_c}{TP_c + FP_c} = \frac{1}{C} \sum_c \text{precision}_c$$

$$\text{macrorecall} = \frac{1}{C} \sum_c \frac{TP_c}{TP_c + FN_c} = \frac{1}{C} \sum_c \text{recall}_c$$

Microaveraging: Collect decisions for all classes, compute contingency table, evaluate.

$$\text{microprecision} = \frac{\sum_c TP_c}{\sum_c TP_c + \sum_c FP_c}$$

$$\text{microrecall} = \frac{\sum_c TP_c}{\sum_c TP_c + \sum_c FN_c}$$

Softmax/MaxEnt

$$p(y | x) = \frac{\exp(\theta_y^T f(x))}{\sum_{y'} \exp(\theta_{y'}^T f(x))}$$

Maximize Log-Likelihood

$$\log \prod_i p_\theta(y_i | x_i) = \sum_i \log p_\theta(y_i | x_i)$$

Inverse of exp
 $\log(\exp(x)) = x$

$$= \sum_i \theta_{y_i}^T f(x_i) - \log Z(x_i)$$

Original maxent equation

$$\frac{\exp(\theta_y^T f(x))}{\sum_{y'} \exp(\theta_{y'}^T f(x))}$$

Differentiating this becomes nicer (even though Z depends on θ)

Minimize Cross Entropy Loss

Cross entropy:
How much \hat{y} differs
from the true y

Model output $\vec{\hat{y}}$ True probability (i.e., correct output) \vec{y}

$$L^{\text{xent}}(\vec{\hat{y}}, \vec{y}) = - \sum_{k=1}^K \vec{y}[k] * \log p(y = k|x)$$

index of "1" indicates correct value $\begin{pmatrix} 0 \\ 0 \\ \dots \\ 1 \\ \dots \\ 0 \end{pmatrix}$ one-hot vector

Probability distribution from model $p(y = k|x)$

Cosine Similarity

$$\cos(x, y) = \frac{\sum_i x_i y_i}{\sqrt{\sum_i x_i^2} \sqrt{\sum_i y_i^2}}$$

	Dim. 1	Dim. 2	Dim. 3
apricot	2	0	0
digital	0	1	2
information	1	6	1

$$\text{cosine}(\text{apricot}, \text{information}) = \frac{2 + 0 + 0}{\sqrt{4 + 0 + 0} \sqrt{1 + 36 + 1}} = 0.1622$$

$$\text{cosine}(\text{digital}, \text{information}) = \frac{0 + 6 + 2}{\sqrt{0 + 1 + 4} \sqrt{1 + 36 + 1}} = 0.5804$$

$$\text{cosine}(\text{apricot}, \text{digital}) = \frac{0 + 0 + 0}{\sqrt{4 + 0 + 0} \sqrt{0 + 1 + 4}} = 0.0$$

Bayes' Rule \rightarrow Naïve Bayes Assumption

Bayes $\hat{c} = \operatorname{argmax}_{c \in C} P(c|d) = \operatorname{argmax}_{c \in C} \frac{P(d|c) \cdot P(c)}{P(d)}$

Naïve Bayes $\hat{c} = \operatorname{argmax}_{c \in C} P(c|d) \approx \operatorname{argmax}_{c \in C} P(d|c) \cdot P(c)$

We can make this assumption because $P(d)$ stays the same regardless of the class!

Perplexity

Lower is better : lower perplexity → less surprised

$$\text{perplexity} = \exp\left(\frac{-1}{M} \sum_{i=1}^M \log p(w_i | h_i)\right)$$

$$= \sqrt[M]{\prod_{i=1}^M \frac{1}{p(w_i | h_i)}}$$

weighted
geometric
average

Trigram MLE

$$\begin{aligned} p(\text{Colorless green ideas sleep furiously}) = & \\ & p(\text{Colorless} \mid \langle \text{BOS} \rangle \langle \text{BOS} \rangle) * \\ & p(\text{green} \mid \langle \text{BOS} \rangle \text{Colorless}) * \\ & p(\text{ideas} \mid \text{Colorless green}) * \\ & p(\text{sleep} \mid \text{green ideas}) * \\ & p(\text{furiously} \mid \text{ideas sleep}) * \\ & p(\langle \text{EOS} \rangle \mid \text{sleep furiously}) \end{aligned}$$

Consistent notation: Pad the left with $\langle \text{BOS} \rangle$ (beginning of sentence) symbols

Fully proper distribution: Pad the right with a single $\langle \text{EOS} \rangle$ symbol

Add- λ estimation

Other names: Laplace
smoothing, Lidstone
smoothing

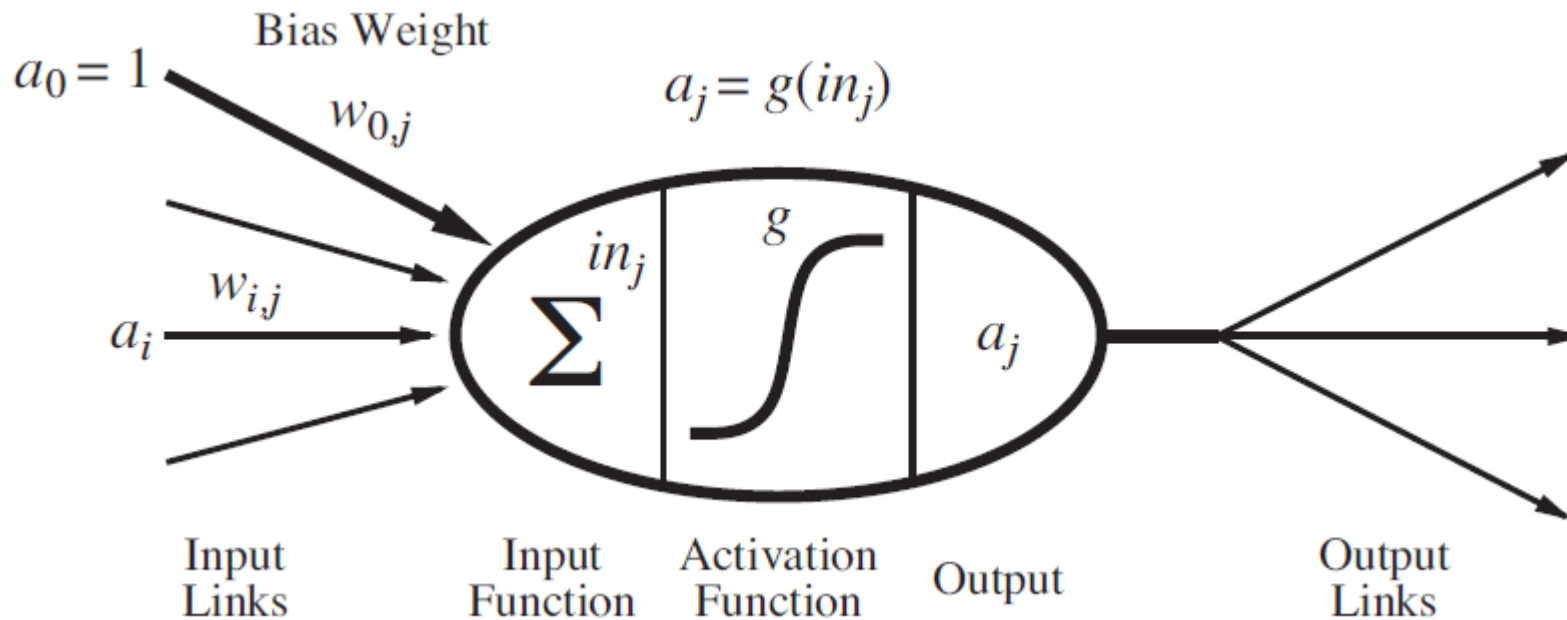
Pretend we saw each
word λ more times than
we did

Add λ to all the counts

$$\begin{aligned} p(\mathbf{z}) &\propto \text{count}(\mathbf{z}) + \lambda \\ &= \frac{\text{count}(\mathbf{z}) + \lambda}{\sum_v (\text{count}(v) + \lambda)} \end{aligned}$$

Biologically-Inspired Learning Models: Neuron Unit

$$in_j = w_{0j} + w_{1j}a_1 + w_{2j}a_2 + \dots + w_{ij}a_i$$



activations
 $0 \leq a_i \leq 1$

weights
 $-\infty < w_{ij} < \infty$

Forward Propagation Example

$$(1) \quad in_j = w_{0j} + w_{1j}a_1 + w_{2j}a_2 + \dots + w_{ij}a_i$$

$$in_{h1} = -.3225$$

$$in_{h2} = -.3075$$

(2) Activation function

Sigmoid

$$g(x) = \frac{1}{1 + e^{-x}}$$

$$\begin{aligned} out_{h1} &= g(in_{h1}) \\ &= \frac{1}{1 + e^{-in_{h1}}} \\ &= \frac{1}{1 + e^{-(-.3225)}} \\ &= .4188 \end{aligned}$$

$$\begin{aligned} out_{h2} &= g(in_{h2}) \\ &= \frac{1}{1 + e^{-in_{h2}}} \\ &= \frac{1}{1 + e^{-(-.3075)}} \\ &= .4237 \end{aligned}$$

For each layer:

1. Calculate the weighted sum of inputs to each neuron unit
2. Evaluate the activation function to determine the output of each neuron unit
3. Use outputs as inputs for the next layer

