

# Naive Bayes' classifier example w/Laplace smoothing

	Cat	Documents
Training	-	just plain boring
	-	entirely predictable and lacks energy
	-	no surprises and very few laughs
	+	very powerful
	+	the most fun film of the summer
Test	?	predictable with no fun

**function** TRAIN NAIVE BAYES(D, C) **returns**  $\log P(c)$  and  $\log P(w|c)$

**for each** class  $c \in C$  # Calculate  $P(c)$  terms

$N_{doc}$  = number of documents in D

$N_c$  = number of documents from D in class c

$\logprior[c] \leftarrow \log \frac{N_c}{N_{doc}}$

$V \leftarrow$  vocabulary of D

$bigdoc[c] \leftarrow$  **append**(d) **for** d  $\in$  D **with** class c

**for each** word w in V # Calculate  $P(w|c)$  terms

$count(w, c) \leftarrow$  # of occurrences of w in  $bigdoc[c]$

$\loglikelihood[w, c] \leftarrow \log \frac{count(w, c) + 1}{\sum_{w' \in V} (count(w', c) + 1)}$

**return**  $\logprior$ ,  $\loglikelihood$ , V

**function** TEST NAIVE BAYES( $testdoc$ ,  $\logprior$ ,  $\loglikelihood$ , C, V) **returns** best c

**for each** class  $c \in C$

$sum[c] \leftarrow \logprior[c]$

**for each** position i in  $testdoc$

$word \leftarrow testdoc[i]$

**if** word  $\in V$

$sum[c] \leftarrow sum[c] + \loglikelihood[word, c]$

**return**  $\operatorname{argmax}_c sum[c]$

# Text classification with discriminative models

- Discriminative modeling approach that learns  $P(y|x)$
- Quick refresher about MLE & discriminative models:

$P(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta})$  where  $\mathbf{x}, \mathbf{y} \sim P_D$ ,  $\mathbf{y} \in \{0,1\}^K \longrightarrow$  one hot encoding

$$\begin{aligned}\hat{\boldsymbol{\theta}} &= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \sum_{i=1}^M -\log P(\mathbf{y}_i | \mathbf{x}_i; \boldsymbol{\theta}) \\ &= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} -\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim P_D} [\log P(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta})] \\ &= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} -\sum_{i=1}^M \mathbf{y} \log P(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta}) \longrightarrow\end{aligned}$$

Cross entropy between data distribution and the model output distribution,  $H(P_D(\mathbf{y} | \mathbf{x}), P(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta}))$ . Note, this holds for arbitrary PMFs (softmax, Bernoulli, etc...). Because the labels,  $\mathbf{y}$ , are one hot, they have zero entropy, and thus in this setting minimizing the cross entropy is equivalent to minimizing the KL divergence,  $D_{KL}(P_D(\mathbf{y} | \mathbf{x}), P(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta}))$ .