























Usi	ng a	test set			
	Set Si	ze Log likelihood			
Training S	et 196	-466.1905			
Test Set	196	-614.6157			
(actually it's a billion qu times less likely)	intillion qu	intillion quintillion	quintillion		
Density estimators can overfit. And the full joint density estimator is the overfittiest of them all!					
Remember the lookup to	able classi	fier with 2^{2^m} hype	otheses?		
ppyright © Andrew W. Moore					



























Joint DE	Naïve DE		
Can model anything	Can model only very boring distributions		
No problem to model	Outside Naïve's scope		
"C is a noisy copy of A"			
Given 100 records and more than 6 Boolean attributes will screw up badly	Given 100 records and 10,000 multivalued attributes will be fine		



































$\begin{aligned} & \text{Getting a posterior probability} \\ & P(Y = v \mid X_1 = u_1 \cdots X_m = u_m) \\ & = \frac{P(X_1 = u_1 \cdots X_m = u_m \mid Y = v) P(Y = v)}{P(X_1 = u_1 \cdots X_m = u_m \mid Y = v) P(Y = v)} \\ & = \frac{P(X_1 = u_1 \cdots X_m = u_m \mid Y = v_j) P(Y = v_j)}{\sum_{j=1}^{n_y} P(X_1 = u_1 \cdots X_m = u_m \mid Y = v_j) P(Y = v_j)} \end{aligned}$

Bayes Classifiers in a nutshell 1. Learn the distribution over inputs for each value Y. 2. This gives $P(X_1, X_2, ..., X_m | Y = v_i)$. 3. Estimate $P(Y = v_i)$. as fraction of records with $Y = v_i$. 4. For a new prediction: $Y_{predict}^{predict} = \operatorname{argmax}_{v} P(Y = v | X_1 = u_1 \cdots X_m = u_m)$ $= \operatorname{argmax}_{v} P(X_1 = u_1 \cdots X_m = u_m | Y = v) P(Y = v)$ v







The attri as 0] "all i ibutes) or 1	oint B rrelevant" s called a,t . v (output	C Results dataset consists b,c,do where a,t c) = 1 with proba	* **All of 40,000 b,c are ge bility 0.75	Ir rec ner 5, 0	relevant" cords and 15 Boolean ated 50-50 randomly with prob 0.25	
Na	ame	Model	Parameters	FracRight			
M	odel1	bayesclass	density=joint submodel=gauss gausstype=general	0.70425	+/-	0.00583537	
							-
Copyright © Andr	rew W. Moo	re					Slide 104

DATE OF A Series Classifier

$$Y^{\text{predict}} = \underset{v}{\operatorname{argmax}} P(X_1 = u_1 \cdots X_m = u_m \mid Y = v) P(Y = v)$$
In the case of the naive Bayes Classifier this can be simplified:

$$Y^{\text{predict}} = \underset{v}{\operatorname{argmax}} P(Y = v) \prod_{j=1}^{n_j} P(X_j = u_j \mid Y = v)$$

Copyright © Andrew W. Moore

Naïve Bayes Classifier $Y^{\text{predict}} = \underset{v}{\operatorname{argmax}} P(X_{1} = u_{1} \cdots X_{m} = u_{m} | Y = v) P(Y = v)$ In the case of the naive Bayes Classifier this can be simplified: $Y^{\text{predict}} = \underset{v}{\operatorname{argmax}} P(Y = v) \prod_{j=1}^{n_{y}} P(X_{j} = u_{j} | Y = v)$ Technical Hint: If you have 10,000 input attributes that product will underflow in floating point math. You should use logs: $Y^{\text{predict}} = \underset{v}{\operatorname{argmax}} \left(\log P(Y = v) + \sum_{j=1}^{n_{y}} \log P(X_{j} = u_{j} | Y = v) \right)$ Substitution of the statement of the statemen

Slide 105













More Facts About Bayes Classifiers

- Many other density estimators can be slotted in*.
- Density estimation can be performed with real-valued inputs*
- Bayes Classifiers can be built with real-valued inputs*
- Rather Technical Complaint: Bayes Classifiers don't try to be maximally discriminative---they merely try to honestly model what's going on*
- Zero probabilities are painful for Joint and Naïve. A hack (justifiable with the magic words "Dirichlet Prior") can help*.
- Naïve Bayes is wonderfully cheap. And survives 10,000 attributes cheerfully!

*See future Andrew Lectures

Slide 113

Copyright © Andrew W. Moore



