

Naïve Bayes

Victor Lavrenko and Nigel Goddard  
School of Informatics

Bayesian classification: components

$$P(y|x) = \frac{P(x|y)P(y)}{\sum_{y'} P(x|y')P(y')}$$

Example:  
y ... patient has Avian flu  
x ... observed symptoms

- **P(y): prior probability of each class**
  - encodes how which classes are common, which are rare
  - apriori much more likely to have common cold than Avian flu
- **P(x|y): class-conditional model**
  - describes how likely to see observation x for class y
  - assuming it's Avian flu, do the symptoms look plausible?
- **P(x): normalize probabilities across observations**
  - does not affect which class is most likely (**arg max**)

Independence assumption

- Compute  $P(x_1 \dots x_n | y)$  for every observation  $x_1 \dots x_n$ 
  - class-conditional "counts", based on training data
  - problem: may not have seen every  $x_1 \dots x_n$  for every y
    - digits:  $2^{400}$  possible black/white patterns (20x20)
    - spam: every possible combination of words:  $2^{10,000}$
  - often have observations for individual  $x_i$  for every class



$$P(x_1 \dots x_d | y) = \underbrace{\prod_{i=1}^d P(x_i | x_1 \dots x_{i-1}, y)}_{\text{chain rule (exact)}} = \underbrace{\prod_{i=1}^d P(x_i | y)}_{\text{independence}}$$

Overview

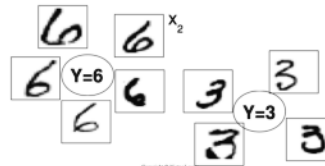
- Naïve Bayes classifier
  - components and their function
  - independence assumption
  - dealing with missing data
- Continuous example
- Discrete example
- Pros and cons

Bayesian classification: normalization

Normalizer:  $P(x) = \sum_{y'} P(x|y')P(y')$

- an "outlier" has a low probability under every class

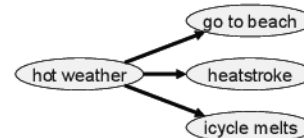
$P(X=x_1 | Y=3) < P(X=x_2 | Y=3)$



normalizer makes  $P(Y=3|X=x_1)$  comparable to non-outliers

Conditional independence

- Probabilities of going to the beach and getting a heat stroke are not independent:  $P(B,S) > P(B)P(S)$
- May be independent if we know the weather is hot  $P(B,S|H) = P(B|H)P(S|H)$
- Hot weather "explains" all the dependence between beach and heatstroke
- In classification:
  - class value explains all the dependence between attributes



Bayesian classification

- Goal: learning function  $f(x) \rightarrow y$ 
  - y ... one of k classes (e.g. spam/ham, digit 0-9)
  - $x = x_1 \dots x_n$  - values of attributes (numeric or categorical)
- Probabilistic classification:
  - most probable class given observation:  $\hat{y} = \arg \max_y P(y|x)$
- Bayesian probability of a class:

$$P(y|x) = \frac{\overbrace{P(x|y)P(y)}^{\text{class model prior}}}{\underbrace{\sum_{y'} P(x|y')P(y')}_{\text{normalizer } P(x)}}$$

Naïve Bayes: a generative model

- A complete probability distribution for each class
  - defines likelihood for any point x
  - P(class) via P(observation)  $P(y|x) \propto P(x|y)P(y)$
  - can "generate" synthetic observations
    - will share many properties of the original data
- Not all probabilistic classifiers do this
  - possible to estimate  $P(y|x)$  directly
  - e.g. logistic regression:

$$P(y|x) = \frac{1}{Z_y} \exp\left(\sum_i \lambda_i g_i(y, x)\right)$$



Overview

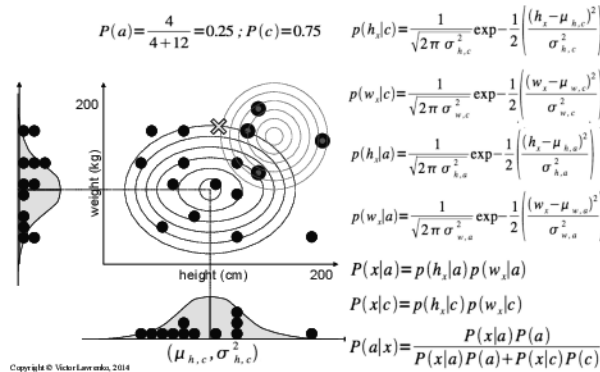
- Naïve Bayes classifier
- Continuous example
  - general concepts
  - working example
  - example of failure
- Discrete example
  - general concepts
  - problems with Naïve Bayes
- Pros and cons

$$P(y|x) = \frac{\overbrace{P(x|y)P(y)}^{\text{class model prior}}}{\underbrace{\sum_{y'} P(x|y')P(y')}_{\text{normalizer } P(x)}}$$

## Continuous example

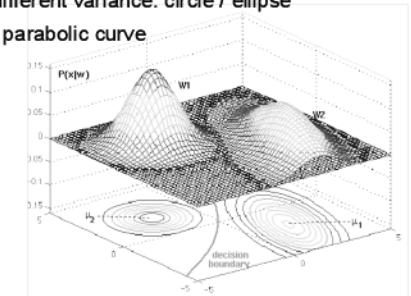
- Distinguish children from adults based on size
  - classes: {a, c}, attributes: height [cm], weight [kg]
  - training examples: {h<sub>p</sub>, w<sub>p</sub>, y<sub>p</sub>}, 4 adults, 12 children
- Class probabilities:  $P(a) = \frac{4}{4+12} = 0.25$ ;  $P(c) = 0.75$
- Model for adults:
  - height ~ Gaussian with mean, variance  $\left\{ \begin{array}{l} \mu_{h,a} = \frac{1}{4} \sum_{i:y=a} h_i \\ \sigma_{h,a}^2 = \frac{1}{4} \sum_{i:y=a} (h_i - \mu_{h,a})^2 \end{array} \right.$
  - weight ~ Gaussian  $(\mu_{w,a}, \sigma_{w,a}^2)$
  - assume height and weight independent
- Model for children: same, using  $(\mu_{h,c}, \sigma_{h,c}^2), (\mu_{w,c}, \sigma_{w,c}^2)$

## Continuous example



## Decision boundary

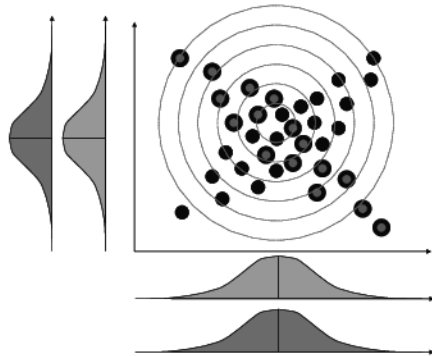
- Different means, same variance: straight line / plane
- Same mean, different variance: circle / ellipse
- General case: parabolic curve



Copyright © Victor Laveenka, 2014

Copyright © Victor Laveenka, 2014

## Problems with Naïve Bayes



Copyright © Victor Laveenka, 2014

## Discrete example: spam

- Separate spam from valid email, attributes = words

D1: "send us your password"	spam	$P(\text{spam}) = 4/6$ $P(\text{ham}) = 2/6$	
D2: "send us your review"	ham	spam	ham
D3: "review your password"	ham	2/4	1/2
D4: "review us"	spam	1/4	2/2
D5: "send your password"	spam	3/4	1/2
D6: "send us your account"	spam	3/4	1/2
		1/4	0/2
		3/4	1/2
		1/4	0/2

$$P(\text{review us}|\text{spam}) = P(0,1,0,1,0,0|\text{spam}) = (1 - \frac{2}{4})(\frac{1}{4})(1 - \frac{3}{4})(\frac{3}{4})(1 - \frac{3}{4})(1 - \frac{1}{4})$$

$$P(\text{review us}|\text{ham}) = P(0,1,0,1,0,0|\text{ham}) = (1 - \frac{1}{2})(\frac{2}{2})(1 - \frac{1}{2})(\frac{1}{2})(1 - \frac{1}{2})(1 - \frac{0}{2})$$

$$P(\text{ham}|\text{review us}) = \frac{0.0625 \times 2/6}{0.0625 \times 2/6 + 0.0044 \times 4/6} = 0.87 \text{ (note identical example)}$$

Copyright © Victor Laveenka, 2014

- Zero-frequency problem

- any mail containing "account" is spam:  $P(\text{account}|\text{ham}) = 0/2$
- solution: never allow zero probabilities
  - Laplace smoothing: add a small positive number to all counts:  $P(w|c) = \frac{\text{num}(w,c) + \epsilon}{\text{num}(c) + 2\epsilon}$
  - may use global statistics in place of  $\epsilon$ :  $\text{num}(w) / \text{num}$
- very common problem (Zipf's law: 50% words occur once)
- Assumes word independence
  - every word contributes independently to  $P(\text{spam}|\text{email})$
  - fooling NB: add lots of "hammy" words into spam email

Copyright © Victor Laveenka, 2014

## Overview

- Naïve Bayes classifier
- Continuous example
- Discrete example
- Pros and Cons
  - dealing with missing data
  - computational cost and incremental updates

## Missing data

- Suppose don't have value for some attribute  $X_i$ 
  - applicant's credit history unknown
  - some medical test not performed on patient
  - how to compute  $P(X_1=x_1 \dots X_j=? \dots X_d=x_d | y)$

- Easy with Naïve Bayes
  - ignore attribute in instance where its value is missing  $P(x_1 \dots \boxed{X_j} \dots x_d | y) = \prod_{i \neq j} P(x_i | y)$
  - compute likelihood based on observed attributes
  - no need to "fill in" or explicitly model missing values
  - based on conditional independence between attributes

Copyright © Victor Laveenka, 2014

Copyright © Victor Laveenka, 2014

## Missing data (2)

- Ex: three coin tosses: Event = { $X_1=H, X_2=?, X_3=T$ }
  - event = head, unknown (either head or tail), tail
  - event = {H,H,T} + {H,T,T}
  - $P(\text{event}) = P(H,H,T) + P(H,T,T)$
- General case:  $X_j$  has missing value

$$\begin{aligned}
 P(x_1 \dots \boxed{X_j} \dots x_d | y) &= P(x_1 | y) \dots P(x_j | y) \dots P(x_d | y) \\
 \sum_{x_j} P(x_1 \dots \boxed{X_j} \dots x_d | y) &= \sum_{x_j} P(x_1 | y) \dots P(x_j | y) \dots P(x_d | y) \\
 &= P(x_1 | y) \dots \left[ \sum_{x_j} P(x_j | y) \right] \dots P(x_d | y) \\
 &= P(x_1 | y) \dots \boxed{1} \dots P(x_d | y)
 \end{aligned}$$

Copyright © Victor Laveenka, 2014

## Summary

- Naïve Bayes classifier

- explicitly handles class priors
- "normalizes" across observations: outliers comparable
- assumption: all dependence is "explained" by class label

$$P(y|x) = \frac{P(x|y)P(y)}{\sum_{y'} P(x|y')P(y')}$$

- Continuous example

- unable to handle correlated data

- Discrete example

- fooled by repetitions
- must deal with zero-frequency problem

- Pros:

- handles missing data
- good computational complexity
- incremental updates

## Computational complexity

- One of the fastest learning methods

- $O(nd+cd)$  training time complexity

- $c$  ... number of classes
- $n$  ... number of instances
- $d$  ... number of dimensions (attributes)
- both learning and prediction
- no hidden constants (number of iterations, etc.)
- testing:  $O(ndc)$

- $O(dc)$  space complexity

- only decision trees are more compact

## Incremental updates

- Allows incremental updates:  $O(d)$  insertion / deletion

- Bernoulli: store raw counts instead of probabilities

- new example of class  $c$ :
  - $n_{cd} += x_d$  for each  $d$  in example,  $n_c += 1$ ,  $n += 1$

- when need to classify:

- $P(x_d \neq 1 | c) = (n_{cd} + \epsilon) / (n_c + 2\epsilon)$
- $P(c) = n_c / n$

- Gaussian: store partial sums instead of mean/variance

- $S_{cd} += X_d$      $S_{cd}^2 += X_d^2$
- when need to classify:  
mean =  $s_{cd} / n$     variance =  $s_{cd}^2 / n - \text{mean}^2$

Copyright © Victor Laveenka, 2014

Copyright © Victor Laveenka, 2014

Copyright © Victor Laveenka, 2014

## General structure for Naïve Bayes

- Task

- $c$ -class classification ( $c \geq 2$ )

- Model structure

- $c \times d$  independent distributions
- continuous: Gaussian, discrete: Bernoulli

- Score function

- class-conditional likelihood

- Optimization / search method

- analytic solution
- Book: section 4.2

Copyright © Victor Laveenka, 2014