

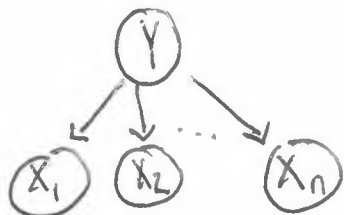
Naive Bayes Classifier

- Family of classification models often used in text classification
- generative
- not Bayesian (but can be made Bayesian)
- major assumption: given the class, y , each feature is conditionally indep. of all other features given y

$$Y \sim \text{cat}(\phi)$$

$$X_j | Y=y \sim P(\Theta_j)$$

Θ_j the j^{th} feature of X



$$y \in \{1, \dots, K\}; \phi \in [0, 1]^K$$

$$X \in \mathbb{R}^n$$

chain rule of probability

$$P(y | X_1, \dots, X_n) = \frac{P(X_1, \dots, X_n | y) P(y)}{P(X_1, \dots, X_n)} = \frac{P(X_1 | y, X_2, \dots) P(X_2 | y, X_3, \dots) \dots P(X_{n-1} | y, X_n) P(X_n | y) P(y)}{P(X_1, \dots, X_n)}$$

$$\stackrel{(\text{c.i.})}{=} \frac{P(y) \prod_{j=1}^n P(X_j | y)}{P(X_1, \dots, X_n)}$$



Class prediction:

$$\hat{y} = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \left\{ P(y; \hat{\phi}) \prod_{j=1}^n P(x_j | y; \hat{\phi}_j) \right\}$$

MLE training (could also use MAP estimate if you want, or go full Bayesian to get a posterior)

let all parameters be in a set, Θ \leftarrow all $x \in \mathbb{R}^n$

$$L(\Theta) = P(\mathcal{D}; \Theta) = P(x^{(1)}, y^{(1)}, x^{(2)}, y^{(2)}, \dots, x^{(n)}, y^{(n)}) \stackrel{\text{iid}}{=} \prod_{i=1}^n P(x^{(i)}, y^{(i)}) = \prod_{i=1}^n P(x^{(i)} | y^{(i)}) P(y^{(i)})$$

$$\stackrel{(NS)}{=} \prod_{i=1}^n \prod_{j=1}^n P(x_j^{(i)} | y^{(i)}) P(y^{(i)}) = \prod_{i=1}^n P(y^{(i)}; \phi) \prod_{j=1}^n P(x_j^{(i)} | y^{(i)}; \phi_j | y^{(i)})$$

\Rightarrow take log and find argmax of Θ , possibly using Lagrange multipliers to make

$$\text{sure } \sum_k \phi_k = 1$$

- MLE estimates are very intuitive and do not require iterative solns. (eg. SGD)

- eg: $\phi_1 = \frac{\text{\# of class 1 in data set}}{\text{\# in data set}}$

- MLEs might benefit from Laplace smoothing, which corresponds to putting a prior on ϕ with ^{not} probability mass $\neq 0$

- depending on application, $P(x|y)$ could be:

- Gaussian
- categorical
- Bernoulli
- etc. ...

eg: email classification

Let $x_j = 0$ denote that the j^{th} word in our dictionary is not in the email and
 $x_j = 1$ denote that it is.

• \therefore a feature vector for an email can be represented as $x = [0, 1, 1, \dots, 0, \dots, 1]^T$ (for example)
 $\in \{0, 1\}^n$

$Y \sim \text{Bern}(\phi)$; $\phi \in [0, 1]$

$x_j | Y=y \sim \text{Bern}(\phi_j | y) \Rightarrow 2 \times (n-1)$ ^{probs. sum to one} parameters to estimate w/ MLE

Pros

- MLE done in closed form, so fast

- Scalable \Rightarrow only requires # parameters \approx # features (i.e., not many parameters to fit)

- means its high bias, low variance

- thus, only need small amount of training data (since not too many params to estimate)

- rather than try to estimate a joint distribution
 $P(x_1, x_2, \dots, x_n, y)$

which is a PDF/PMF in a very high dim. space (which would take many parameters to actually fit and lots of data), the "Naive" assumption allows us to decouple the problem to estimating n 1-D (univariate) PDFs/PMFs.

$$P(x_1, y) P(x_2, y) \dots P(x_n, y)$$

Cons

- many times cond. independence does not hold (but NB still seems to perform well in many applications)
- high ~~bias~~ bias model (sometimes too simple) so that other, more complicated models might perform better (e.g. RF)