

03/02/15

1

BAYES CLASSIFIERS + NAIVE BAYES

① Bayes Error

→ The best ANY ML algorithm can do.

Recall Expected Loss/Error

$\vec{x}, y \sim P^*(x, y)$ [Unknown] (Using P^* just to emphasize that this is reality, not from model)

$$E[\text{Loss}] = E_{P^*(\vec{x}, y)} [L(y, \hat{y})]$$

reality
prediction

Bayes Error
≡ min possible value of this

$$= \int_{\vec{x}} \int_y L(y, \hat{y}) P^*(\vec{x}, y) d\vec{x} dy$$

Space of \vec{x}
Say \mathbb{R}^d
Space of Y

→ Assume Classification so $Y = \{1, 2, \dots, k\}$
[Replace \int_y with $\sum_{y=1}^k$]

→ Assume 0-1 Loss i.e. $L(y, \hat{y}) = \begin{cases} 1 & y \neq \hat{y} \\ 0 & \text{else} \end{cases}$

$$\Rightarrow E[\text{Loss}] = \int_{\vec{x}} \sum_{y=1}^k L(y, \hat{y}) P^*(\vec{x}, y) d\vec{x}$$

[Mix of PMF & PDF]

$$E[\text{Loss}] = \int_{\vec{x}} \left[\sum_{y=1}^k L(y, \hat{y}) p^*(y|\vec{x}) \right] p(\vec{x}) d\vec{x}$$

looks like $E_{p(\vec{x})}[f(\vec{x})]$

$$= E_{p^*(\vec{x})} \left[\sum_{y=1}^k L(y, \hat{y}) p^*(y|\vec{x}) \right]$$

And what happens when I predict \hat{y}

For this \vec{x} how likely is the output to be y (in reality)

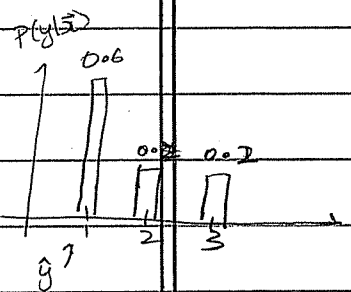
$$= E_{p^*(\vec{x})} \left[L(1, \hat{y}) p^*(1|\vec{x}) + L(2, \hat{y}) p^*(2|\vec{x}) + \dots + L(k, \hat{y}) p^*(k|\vec{x}) \right]$$

either 0 or 1

$$= E_{p^*(\vec{x})} \left[1 - p^*(Y = \hat{y} | \vec{x}) \right]$$

So $E[\text{Loss}] = E_{p^*(\vec{x})} \left[1 - p^*(Y = \hat{y} | \vec{x}) \right]$

$$\text{So } \min_{\hat{y}} E[\text{Loss}] = \max_{\hat{y}} p^*(Y = \hat{y} | \vec{x})$$



Bayes Error
[Error of the best possible prediction]

MAP [Maximum-a-posteriori]
aka "pick the most confident class"

\Rightarrow So \hat{y}_{MAP} is Bayes Optimal!

→ Bayes Classifier: the classifier that achieves Bayes Error.

[Similarly Bayes Regressor]

→ Note: \hat{y}_{MAP} is Bayes Optimal FOR 0/1-Loss.

Different Loss \Rightarrow Different Bayes Classifiers

→ Great theory! How do I implement this?

Problem: Don't know $P^*(Y|\vec{X}=\vec{x})$

Solution: → Generative Approach
estimate

$P(\vec{x}|y)$ from \Rightarrow Bayes $\Rightarrow P(y|\vec{x}) \Rightarrow \hat{y}_{MAP}$
 $P(y)$ data Rule

→ Discriminative Approach:

estimate

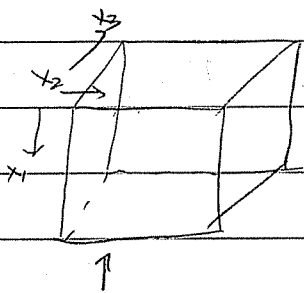
$P(y|\vec{x}) \Rightarrow$ predict & done
from data

② Bayes Classifiers are hard to learn: A counting argument

Say $\vec{X} = \vec{x} \in \{0,1\}^d$ d -binary-features
 $Y = y \in \{1,2,\dots,k\}$ k classes

Need to estimate: $P(Y=y)$ & $P(\vec{X}=\vec{x} | Y=y)$

How many parameters? $(k-1) + (2^d - 1) \cdot k$



Very sparse table

If $d=100$ $2^{100} \approx 10^{30} \gg \gg$ data available

lots of parameters \Rightarrow high variance to overfitting

③ Naive Bayes to the Rescue.

Assume $P(x_1=x_1, \dots, x_d=x_d | Y=y) = \prod_{j=1}^d P(x_j=x_j | Y=y) \forall x_j, y$

"Features are conditionally indep given class"

If I tell you $x_2=x_2$ you don't find out anything about $x_1=x_1$ any more than $Y=y$ already told you

Note: this does not mean $P(x_1 | x_2) = P(x_1)$

How many parameters in NB?

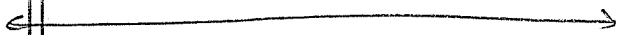
$$d \cdot \underbrace{(2^d - 1) \cdot k}_{P(x_i=x_i | Y=c)} + \underbrace{(k-1)}_{P(Y=y)}$$

Decision Rule:

$$\begin{aligned} \hat{y}_{MAP} &= \operatorname{argmax}_y P(Y=y | \bar{X}=\bar{x}) \\ &= \operatorname{argmax}_y \left[\frac{P(\bar{X}=\bar{x} | Y=y) P(Y=y)}{P(\bar{X}=\bar{x})} \right] \\ &= \operatorname{argmax}_y \prod_{j=1}^d P(X_j=x_j | Y=y) \cdot P(Y=y) \end{aligned}$$

(sometimes also written as)

$$= \operatorname{argmax}_y \left[\sum_{j=1}^d \log P(X_j=x_j | Y=y) + \log P(Y=y) \right]$$



④ MLE/MAP for estimating NB parameters

→ Given dataset $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$

→ What are the parameters?

→ $P(Y)$ ← Categorical $\begin{bmatrix} 1 \\ \vdots \\ P(Y=c) \\ \vdots \\ K \end{bmatrix}$ ← π_c (shorthand)

→ $P(X_j=a | Y=c)$

	Y			
	1	2	...	k
X_j	1			
2				
$ X_j $				

$= \theta_a^{jc}$

shorthand

$$\sum_{a=1}^{|X_j|} \theta_a^{jc} = 1$$

↑ Each column is categorical

All parameters $\Theta = \{\pi_1, \dots, \pi_k, \theta_a^{ic} \forall j, c, a\}$

We know how to estimate categorical parameters
HW1 Q2.

$\hat{\Theta}$
MLE

$$\hat{\pi}_c = P(Y=c) = \frac{\text{Count}(Y=c)}{n}$$

add pseudo-counts for $\hat{\Theta}_{MAP}$

$$\hat{\theta}_a^{ic} = P(X_j=a | Y=c) = \frac{\text{Count}(X_j=a, Y=c)}{\text{Count}(Y=c)}$$

←

⑤ Bag-of-words model

NB assumes $X_i \perp X_j | Y$

ie $P(X_i | X_j, Y) = P(X_i | Y)$

BoW additionally assumes $P(X_i | Y) = P(X_j | Y)$

that all dimensions are essentially the same. More sharing. Less Parameters.

Less expressive model class. Less variance.

Less overfitting.