

Bayes optimal classifier

Naïve Bayes

What's learning, revisited

Machine Learning – 10701/15781
Carlos Guestrin
Carnegie Mellon University
September 21st, 2009

©Carlos Guestrin 2005-2009

1

Classification

- Learn: $h: X \mapsto Y$
 - X – features
 - Y – target classes

GPA
Grade 10701

instead of continuous
as in regression
↳ discrete
↳ hired, not hired!

- Suppose you know $P(Y|X)$ exactly, how should you classify?

- Bayes classifier:

$$y^* = \underset{y}{\operatorname{argmax}} P(Y=y | X=x)$$

$$P(Y=\text{hired} | \text{GPA}=4.0, \text{10701}=8) = 0.01$$

- Why?

Optimal classification

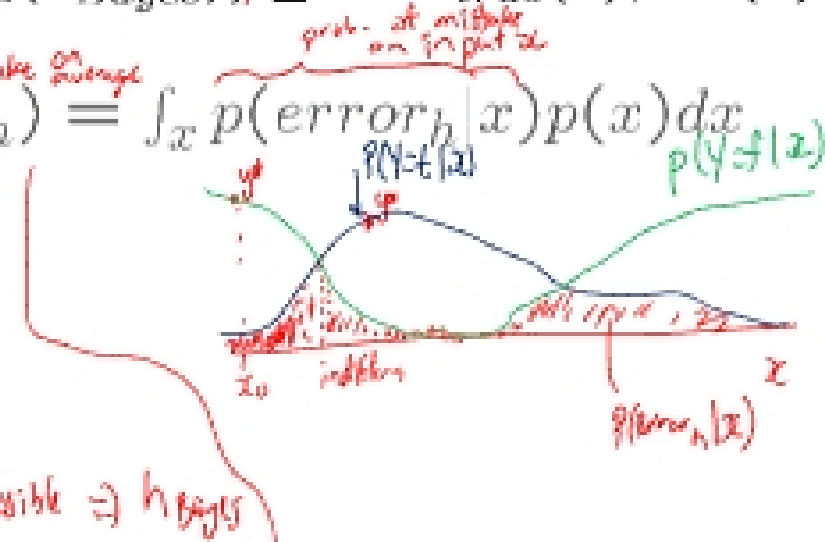
- **Theorem:** Bayes classifier h_{Bayes} is optimal!

$$h_{\text{Bayes}}(x) = \underset{y}{\text{argmax}} P(Y=y | X=x)$$

- That is $\text{error}_{\text{true}}(h_{\text{Bayes}}) \leq \text{error}_{\text{true}}(h), \forall h(x)$

- **Proof:** $p(\text{error}_h) = \int_x p(\text{error}_h | x) p(x) dx$

if I say $h(x) = \text{hired}$
 $p(\text{error}_h | x) \leftarrow$ small
 $= P(Y = \text{not hired} | x)$
 as small as possible $\Rightarrow h_{\text{Bayes}}$



Bayes Rule

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

\nearrow we want this

likelihood prior
 $P(\text{CR} = \text{A.O.} | \text{CR} = \text{hired})$
 $P(\text{hired})$

Which is shorthand for:

$$(\forall i, j) P(Y = y_i | X = x_j) = \frac{P(X = x_j | Y = y_i) P(Y = y_i)}{P(X = x_j)}$$

How hard is it to learn the optimal classifier?

■ Data =

| Sky | Temp | Humid | Wind | Water | Forest | EnjoySport |
|-------|------|--------|--------|-------|--------|------------|
| Sunny | Warm | Normal | Strong | Warm | Some | Yes |
| Sunny | Warm | High | Strong | Warm | Some | Yes |
| Rainy | Cold | High | Strong | Warm | Change | No |
| Sunny | Warm | High | Strong | Cool | Change | Yes |

■ How do we represent these? How many parameters?

(not k , □ Prior, $P(Y)$:

because $\sum_y P(y) = 1$ □ Suppose Y is composed of k classes

params: $\Rightarrow P(Y=y) \in$ prob of outcome y in general
 $k-1$

□ Likelihood, $P(X|Y)$:

■ Suppose X is composed of n binary features

$$P(\text{Sky}=s, \text{Temp}=w, \text{Humid}=N, \text{Wind}=S, \text{Water}=w, \text{Forest}=S | E=\text{yes})$$

$$P(x_1=x_1, x_2=x_2, \dots, x_n=x_n | Y=y) \leftarrow (2^n - 1)k$$

2^n settings

param
a lot,
a lot

■ Complex model ! High variance with limited data!!!

Conditional Independence

indep. $X \perp Y$
 $P(X|Y) = P(X)$

Symmetric relation

■ X is **conditionally Independent** of Y given Z , if the probability distribution governing X is independent of the value of Y , given the value of Z
 $(\forall i, j, k) P(X = i | Y = j, Z = k) = P(X = i | Z = k)$

■ e.g., $P(\text{Thunder} | \text{Rain}, \text{Lightning}) = P(\text{Thunder} | \text{Lightning})$

Thunder indep. of Rain given Lightning

not Thunder indep. Rain.

■ Equivalent to:

$$P(X, Y | Z) = P(X | Z)P(Y | Z)$$