# A Probabilistic View of Machine Learning (1/2)

**CMSC 422** 

MARINE CARPUAT

marine@cs.umd.edu

# Today's topics

- Bayes rule review
- A probabilistic view of machine learning
  - Joint Distributions
  - Bayes optimal classifier
- Statistical Estimation
  - Maximum likelihood estimates
  - Derive relative frequency as the solution to a constrained optimization problem

# Bayes Rule

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$
 Bayes' rule

we call P(A) the "prior"

and P(A|B) the "posterior"



Bayes, Thomas (1763) An essay towards solving a problem in the doctrine of chances. *Philosophical Transactions of the Royal Society of London*, **53:370-418** 

...by no means merely a curious speculation in the doctrine of chances, but necessary to be solved in order to a sure foundation for all our reasonings concerning past facts, and what is likely to be hereafter.... necessary to be considered by any that would give a clear account of the strength of *analogical* or *inductive reasoning...* 

# Exercise: Applying Bayes Rule

Consider the 2 random variables

A = 1 if you have the flu, 0 otherwise

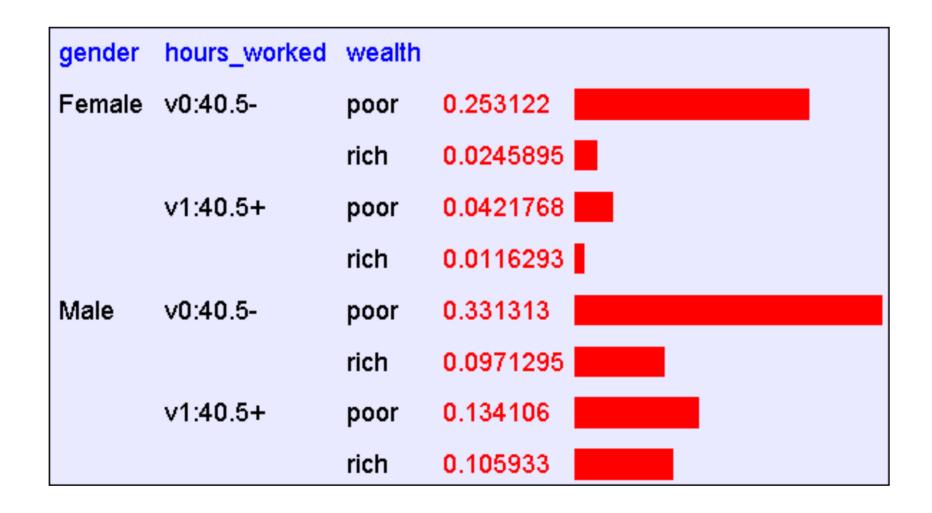
B = 1 if you just coughed, 0 otherwise

Assume

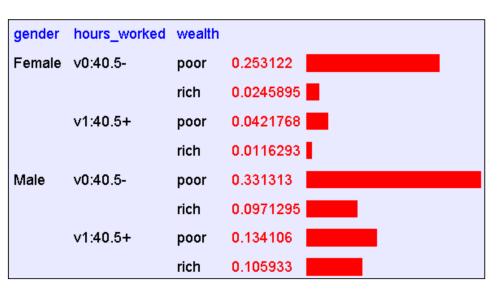
$$P(A = 1) = 0.05$$
  
 $P(B = 1|A = 1) = 0.8$   
 $P(B = 1|A = 0) = 0.2$ 

• What is P(A = 1|B = 1)?

# Using a Joint Distribution



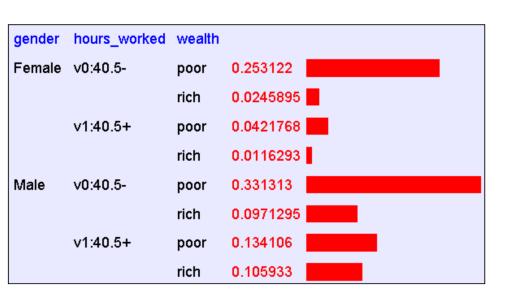
# Using a Joint Distribution



Given the joint
 distribution, we can find
 the probability of any
 logical expression E
 involving these variables

$$P(E) = \sum_{\text{rows matching } E} P(\text{row})$$

### Using a Joint Distribution



Given the joint distribution, we can make inferences

- E.g., P(Male|Poor)?
- Or P(Wealth | Gender, Hours)?

# Recall: Machine Learning as Function Approximation

#### Problem setting

- Set of possible instances X
- Unknown target function  $f: X \to Y$
- Set of function hypotheses  $H = \{h \mid h: X \rightarrow Y\}$

#### Input

• Training examples  $\{(x^{(1)},y^{(1)}),...(x^{(N)},y^{(N)})\}$  of unknown target function f

#### Output

• Hypothesis  $h \in H$  that best approximates target function f

# Recall: Formal Definition of Binary Classification (from CIML)

#### TASK: BINARY CLASSIFICATION

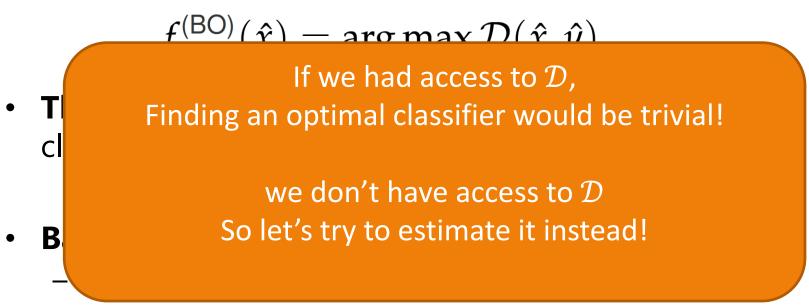
#### Given:

- 1. An input space  $\mathcal{X}$
- 2. An unknown distribution  $\mathcal{D}$  over  $\mathcal{X} \times \{-1, +1\}$

*Compute:* A function f minimizing:  $\mathbb{E}_{(x,y)\sim\mathcal{D}}[f(x)\neq y]$ 

# The Bayes Optimal Classifier

- Assume we know the data generating distribution  ${\cal D}$
- We define the Bayes Optimal classifier as



Best error rate we can ever hope to achieve under zero/one loss

# What does "training" mean in probabilistic settings?

- Training = estimating  $\mathcal{D}$  from a finite training set
  - We typically assume that  $\ensuremath{\mathcal{D}}$  comes from a specific family of probability distributions
    - e.g., Bernouilli, Gaussian, etc
  - Learning means inferring parameters of that distributions
    - e.g., mean and covariance of the Gaussian

# Training assumption: training examples are iid

# Independently and Identically distributed

– i.e. as we draw a sequence of examples from  $\mathcal{D}$ , the n-th draw is independent from the previous n-1 sample

- This assumption is usually false!
  - But sufficiently close to true to be useful

How can we estimate the joint probability distribution from data?

What are the challenges?

### Maximum Likelihood Estimation

Find the parameters that maximize the probability of the data

### Maximum Likelihood Estimates



X=1 X=0  $P(X=1) = \theta$   $P(X=0) = 1-\theta$ (Bernoulli) Each coin flip yields a Boolean value for X

$$X \sim \text{Bernouilli: } P(X) = \theta^X (1 - \theta)^X$$

Given a data set D of iid flips, which contains  $\alpha_1$  ones and  $\alpha_0$  zeros  $P_{\theta}(D) = \theta^{\alpha_1} (1 - \theta)^{\alpha_0}$ 

$$\hat{\theta}_{MLE} = argmax_{\theta} P_{\theta}(D) = \frac{\alpha_1}{\alpha_1 + \alpha_0}$$

### Maximum Likelihood Estimation

Example: how to model a k-sided die?
 (on board)

# Today's topics

- Bayes rule review
- A probabilistic view of machine learning
  - Joint Distributions
  - Bayes optimal classifier
- Statistical Estimation
  - Maximum likelihood estimates
  - Derive relative frequency as the solution to a constrained optimization problem