# Machine Learning 10-601

Tom M. Mitchell
Machine Learning Department
Carnegie Mellon University

January 26, 2015

Today:
- Bayes Classifiers
- Conditional Independence
- Naïve Bayes

Readings:

Mitchell:
"Naïve Bayes and Logistic Regression"
(available on class website)

# Two Principles for Estimating Parameters

- Maximum Likelihood Estimate (MLE): choose $\theta$ that maximizes probability of observed data $\mathcal{D}$

$$\frac{\alpha_1}{\alpha_1 + \alpha_2} \qquad \hat{\theta} = \arg\max_\theta P(\mathcal{D} \mid \theta) \qquad P(Data \mid Model)$$

- Maximum a Posteriori (MAP) estimate: choose $\theta$ that is most probable given prior probability and the data

$$\hat{\theta} = \arg\max_\theta P(\theta \mid \mathcal{D}) \qquad P(Model \mid Data)$$

$$\frac{\alpha_1 + \beta_1}{(\alpha_1 + \beta_1) + (\alpha_2 + \beta_2)} \qquad = \arg\max_\theta = \frac{P(\mathcal{D} \mid \theta)P(\theta)}{P(\mathcal{D})}$$

# Maximum Likelihood Estimate

X=1     X=0

P(X=1) = θ

P(X=0) = 1-θ

(Bernoulli)

- Each flip yields boolean value for $X$

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

- Data set $D$ of independent, identically distributed (iid) flips produces $\alpha_1$ ones, $\alpha_0$ zeros

$$P(D|\theta) = P(\alpha_1, \alpha_0 | \theta) = \theta^{\alpha_1}(1-\theta)^{\alpha_0}$$

$$\hat{\theta}^{MLE} = \arg\max_\theta P(D|\theta) = \frac{\alpha_1}{\alpha_1 + \alpha_0}$$

# Maximum A Posteriori (MAP) Estimate

X=1    X=0

- Data set $D$ of independent, identically distributed (iid) flips produces $\alpha_1$ ones, $\alpha_0$ zeros

$$P(D|\theta) = P(\alpha_1, \alpha_0|\theta) = \theta^{\alpha_1}(1-\theta)^{\alpha_0}$$

- Assume prior $P(\theta) = Beta(\beta_1, \beta_0) = \frac{1}{B(\beta_1,\beta_0)} \theta^{\beta_1-1}(1-\theta)^{\beta_0-1}$

- Then

$$\hat{\theta}^{MAP} = \arg\max_{\theta} P(D|\theta)P(\theta) = \frac{\alpha_1 + \beta_1 - 1}{(\alpha_1 + \beta_1 - 1) + (\alpha_0 + \beta_0 - 1)}$$

(handwritten annotation: $\theta^{\alpha_1 + \beta_1 - 1} [(1-\theta)]^{\alpha_0 + \beta_0 - 1}$)

(like MLE, but hallucinating $\beta_1-1$ additional heads, $\beta_0-1$ additional tails)

# Let's learn classifiers by learning P(Y|X)

Consider Y=Wealth,  X=<Gender, HoursWorked>   $P(W | G, H)$

| gender | hours_worked | wealth | | |
|--------|-------------|--------|--------|---|
| Female | v0:40.5- | poor | 0.253122 | |
| | | rich | 0.0245895 | |
| | v1:40.5+ | poor | 0.0421768 | |
| | | rich | 0.0116293 | |
| Male | v0:40.5- | poor | 0.331313 | |
| | | rich | 0.0971295 | |
| | v1:40.5+ | poor | 0.134106 | |
| | | rich | 0.105933 | |

$2^3 = 8$

#params = 7

| Gender | HrsWorked | P(rich \| G,HW) | P(poor \| G,HW) |
|--------|-----------|----------------|----------------|
| F | <40.5 | .09 | .91 |
| F | >40.5 | .21 | .79 |
| M | <40.5 | .23 | .77 |
| M | >40.5 | .38 | .62 |

$1 - P$

#params = 4

# How many parameters must we estimate?

| Gender | HrsWorked | P(rich \| G,HW) | P(poor \| G,HW) |
|--------|-----------|-----------------|------------------|
| F | <40.5 | .09 | .91 |
| F | >40.5 | .21 | .79 |
| M | <40.5 | .23 | .77 |
| M | >40.5 | .38 | .62 |

Suppose $X = <X_1, \ldots X_n>$

where $X_i$ and $Y$ are boolean <u>RV</u>'s

$$G_T(x) = \text{argmin}_{Y \in G} P(Y \mid x)$$

To estimate $P(Y \mid X_1, X_2, \ldots X_n)$

$$Y = 1 : \quad 2^n$$
$$Y = 0 : \quad 2^n$$

$$\# params = 2^n$$

If we have 30 boolean $X_i$'s: $P(Y \mid X_1, X_2, \ldots X_{30})$

$$2^{30} = (2^{10})^3 \simeq 10^9$$

# Bayes Rule

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

Which is shorthand for:

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

Equivalently:

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

# Can we reduce params using Bayes Rule?

Suppose $X = <X_1, \ldots X_n>$
where $X_i$ and $Y$ are boolean RV's

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

$2^n$

$2^{n+1} - 1$

How many parameters to define $P(X_1, \ldots X_n | Y)$?

$Y = 1: \quad 2^n - 1$

$Y = 0: \quad 2^n - 1 \quad \rightarrow \quad 2 \cdot (2^n - 1)$

How many parameters to define $P(Y)$?

$1$

# Naïve Bayes

Naïve Bayes assumes

$$P(X_1 \ldots X_n | Y) = \prod_i P(X_i | Y)$$

i.e., that $X_i$ and $X_j$ are conditionally independent given Y, for all $i \neq j$

# Conditional Independence

Definition: X is <u>conditionally independent</u> of Y given Z, if
the probability distribution governing X is independent
of the value of Y, given the value of Z

$$(\forall i, j, k) P(X = x_i | Y = y_j, Z = z_k) = P(X = x_i | Z = z_k)$$

Which we often write

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

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

$$P(Y|Z)$$

E.g.,

$$P(Thunder | Rain, Lightning) = P(Thunder | Lightning)$$

$$P(T, R | L) = P(T|L) P(R|L)$$

Naïve Bayes uses assumption that the $X_i$ are conditionally independent, given Y.   E.g., $P(X_1|X_2, Y) = P(X_1|Y)$

Given this assumption, then:

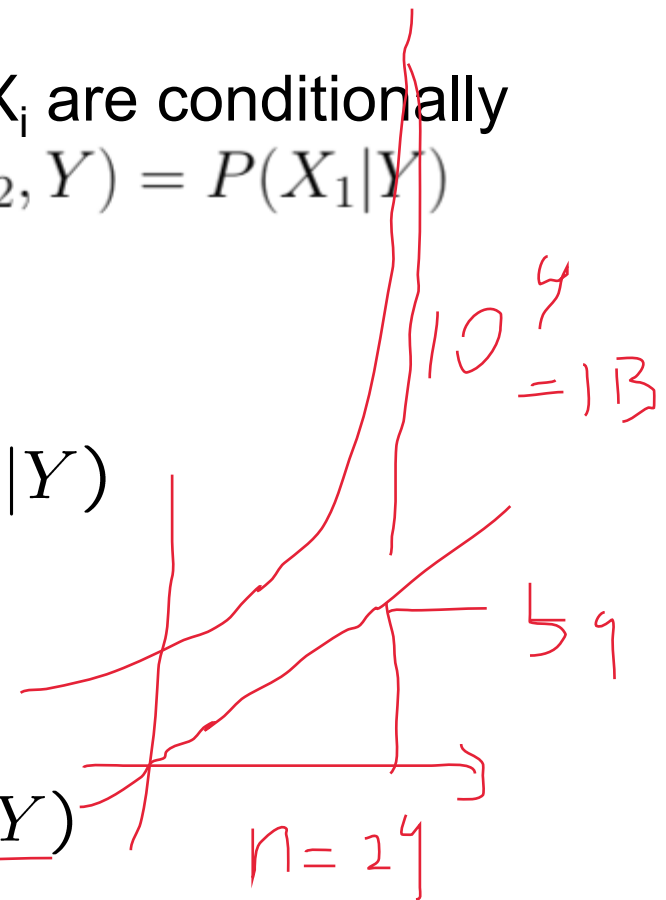$P(X_1, X_2|Y) = P(X_1|X_2, Y) P(X_2|Y)$

$= P(X_1|Y) P(X_2|Y)$

$P(X_1 \ldots, X_n|Y) = \prod_{i=1}^{n} P(X_i|Y)$

Naïve Bayes uses assumption that the $X_i$ are conditionally independent, given Y.  E.g., $P(X_1|X_2, Y) = P(X_1|Y)$

Given this assumption, then:

$$P(X_1, X_2|Y) = P(X_1|X_2, Y)P(X_2|Y)$$
$$= P(X_1|Y)P(X_2|Y)$$

in general:  $P(X_1...X_n|Y) = \prod_i P(X_i|Y)$

Naïve Bayes uses assumption that the $X_i$ are conditionally independent, given Y.  E.g., $P(X_1|X_2, Y) = P(X_1|Y)$

Given this assumption, then:

$$P(X_1, X_2|Y) = P(X_1|X_2, Y)P(X_2|Y)$$
$$= P(X_1|Y)P(X_2|Y)$$

in general:  $P(X_1...X_n|Y) = \prod_i P(X_i|Y)$

$10\frac{y}{} = 13$

$5\,9$

$n = 24$

How many parameters to describe $P(X_1...X_n|Y)$?  $P(Y)$?
- Without conditional indep assumption? $2(2^n - 1) \quad + \quad 1$
- With conditional indep assumption? $2n \quad + \quad 1$

$Y = 1: \quad 2n/2 = r$

$Y = a: \quad 2n/2 = r$

# Naïve Bayes in a Nutshell

Bayes rule:

$$P(Y = y_k | X_1 \dots X_n) = \frac{P(Y = y_k)P(X_1 \dots X_n | Y = y_k)}{\sum_j P(Y = y_j)P(X_1 \dots X_n | Y = y_j)}$$

Assuming conditional independence among $X_i$'s:

$$P(Y = y_k | X_1 \dots X_n) = \frac{P(Y = y_k)\prod_i P(X_i | Y = y_k)}{\sum_j P(Y = y_j)\prod_i P(X_i | Y = y_j)}$$

So, to pick most probable Y for $X^{new} = <X_1, \dots, X_n>$

$$Y^{new} \leftarrow \arg\max_{y_k} P(Y = y_k)\prod_i P(X_i^{new} | Y = y_k)$$

$$P(X_i^{new} = x_{ij} | Y = y_k)$$

# Naïve Bayes Algorithm – discrete $X_i$

$$X_i \in \{1, \ldots, J\}, \quad Y \in \{1 \ldots K\}$$

- Train Naïve Bayes (examples)

  for each* value $y_k$

   estimate $\pi_k \equiv P(Y = y_k)$ · $K$

   for each* value $x_{ij}$ of each attribute $X_i$

   estimate $\theta_{ijk} \equiv P(X_i = x_{ij} | Y = y_k)$ $n J K$

- Classify ($X^{new}$)

$$Y^{new} \leftarrow \arg\max_{y_k} \; P(Y = y_k) \prod_i P(X_i^{new} | Y = y_k)$$

$$Y^{new} \leftarrow \arg\max_{y_k} \; \pi_k \prod_i \theta_{ijk}$$

* probabilities must sum to 1, so need estimate only n-1 of these...

# Estimating Parameters: $Y, X_i$ discrete-valued

Maximum likelihood estimates (MLE's):

$$\hat{\pi}_k = \hat{P}(Y = y_k) = \frac{\#D\{Y = y_k\}}{|D|}$$

$$\hat{\theta}_{ijk} = \hat{P}(X_i = x_{ij} | Y = y_k) = \frac{\#D\{X_i = x_{ij} \wedge Y = y_k\}}{\#D\{Y = y_k\}}$$

Number of items in dataset D for which Y=y$_k$

# Naïve Bayes: Subtlety #1

Often the $X_i$ are not really conditionally independent

- We use Naïve Bayes in many cases anyway, and it often works pretty well
  - often the right classification, even when not the right probability (see [Domingos&Pazzani, 1996])

- What is effect on estimated P(Y|X)?
  - Extreme case: what if we add two copies: $X_i = X_k$

$$P(Y=1 \mid X) \propto P(Y=1) \, P(X_i \mid Y=1) \, P(X_2 \mid Y=1) \cdots$$

# Naïve Bayes: Subtlety #2

If unlucky, our MLE estimate for $P(X_i | Y)$ might be zero.
(for example, $X_i$ = *birthdate*.  $X_i$ = *Jan_25_1992*)

$$\longrightarrow P(X_i = - | Y) = 0 = \frac{1}{365}$$

- Why worry about just one parameter out of many?

$$P(Y=1 | X_1 \ldots X_n) = \frac{P(Y) \prod_i P(X_i | Y)}{P(X)} = 0$$

- What can be done to address this?  ① High-dim

$$MLE \longleftarrow prior$$
$$\Downarrow$$
$$MAP$$

② Spare

# Estimating Parameters

- Maximum Likelihood Estimate (MLE): choose θ that maximizes probability of observed data $\mathcal{D}$

$$\widehat{\theta} = \arg\max_{\theta} P(\mathcal{D} \mid \theta)$$

- Maximum a Posteriori (MAP) estimate: choose θ that is most probable given prior probability and the data

$$\widehat{\theta} = \arg\max_{\theta} P(\theta \mid \mathcal{D})$$

$$= \arg\max_{\theta} = \frac{P(\mathcal{D} \mid \theta)P(\theta)}{P(\mathcal{D})}$$

# Estimating Parameters: $Y, X_i$ discrete-valued

## Maximum likelihood estimates:

$$\hat{\pi}_k = \hat{P}(Y = y_k) = \frac{\#D\{Y = y_k\}}{|D|}$$

$$\hat{\theta}_{ijk} = \hat{P}(X_i = x_j | Y = y_k) = \frac{\#D\{X_i = x_j \wedge Y = y_k\}}{\#D\{Y = y_k\}}$$

## MAP estimates (Beta, Dirichlet priors):

$$\hat{\pi}_k = \hat{P}(Y = y_k) = \frac{\#D\{Y = y_k\} + (\beta_k - 1)}{|D| + \sum_m (\beta_m - 1)}$$

Only difference: "imaginary" examples

$$\hat{\theta}_{ijk} = \hat{P}(X_i = x_j | Y = y_k) = \frac{\#D\{X_i = x_j \wedge Y = y_k\} + (\beta_k - 1)}{\#D\{Y = y_k\} + \sum_m (\beta_m - 1)}$$

# What you should know:

- Training and using classifiers based on Bayes rule

- Conditional independence
  - What it is
  - Why it's important

- Naïve Bayes
  - What it is
  - Why we use it so much
  - Training using MLE, MAP estimates
  - Discrete variables and continuous (Gaussian)

$$\left( 2\left( 2^n - 1 \right) + 1 \right)$$

$$2n + 1$$

$$\hat{G}(x) = \arg\max P(Y|X) \propto P(Y) P(X|Y)$$

$$\prod P(X_i|Y)$$

MLE   $\theta$  $\bar{\pi}$   GNB
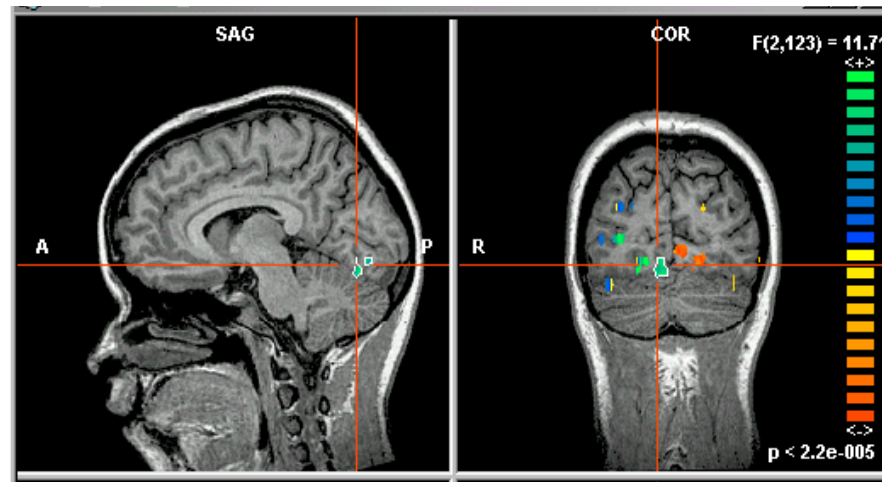
MAP   $\theta$  $\pi$

# Questions:

- How can we extend Naïve Bayes if just 2 of the $X_i$'s are <u>dependent</u>?

- What does the decision surface of a Naïve Bayes classifier look like?

- What error will the classifier achieve if Naïve Bayes assumption is satisfied and we have infinite training data?

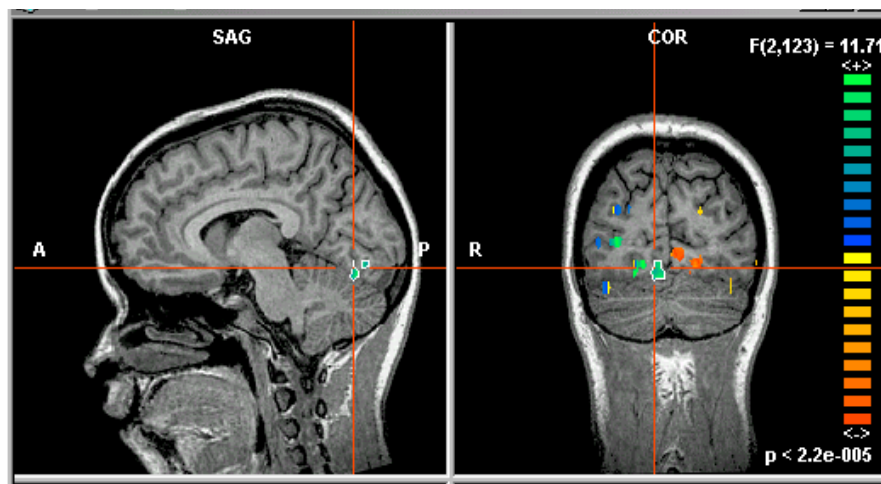- Can you use Naïve Bayes for a combination of discrete and real-valued $X_i$?

# What if we have continuous $X_i$?

Eg., image classification: $X_i$ is i<sup>th</sup> pixel

# What if we have continuous $X_i$ ?

image classification: $X_i$ is i[th] pixel, Y = mental state



Still have:

$$P(Y = y_k | X_1 \ldots X_n) = \frac{P(Y = y_k) \prod_i P(X_i | Y = y_k)}{\sum_j P(Y = y_j) \prod_i P(X_i | Y = y_j)}$$

Just need to decide how to represent P($X_i$ | Y)

# What if we have continuous $X_i$?

Eg., image classification: $X_i$ is i[th] pixel

Gaussian Naïve Bayes (GNB): assume

$$P(X_i = x \mid Y = y_k) = \frac{1}{\sigma_{ik}\sqrt{2\pi}} \; e^{\frac{-(x-\mu_{ik})^2}{2\sigma_{ik}^2}}$$

Sometimes assume $\sigma_{ik}$
- is independent of Y (i.e., $\sigma_i$),
- or independent of $X_i$ (i.e., $\sigma_k$)
- or both (i.e., $\sigma$)

# Gaussian Naïve Bayes Algorithm – continuous $X_i$
(but still discrete Y)

- Train Naïve Bayes (examples)

  for each value $y_k$

  estimate* $\pi_k \equiv P(Y = y_k)$

  for each attribute $X_i$ estimate

  class conditional mean $\mu_{ik}$ , variance $\sigma_{ik}$

- Classify ($X^{new}$)

$$Y^{new} \leftarrow \arg\max_{y_k} \ P(Y = y_k) \prod_i P(X_i^{new} | Y = y_k)$$

$$Y^{new} \leftarrow \arg\max_{y_k} \ \pi_k \prod_i Normal(X_i^{new}, \mu_{ik}, \sigma_{ik})$$

* probabilities must sum to 1, so need estimate only n-1 parameters...

# Estimating Parameters: $Y$ discrete, $X_i$ continuous

Maximum likelihood estimates:

jth training example

$$\widehat{\mu}_{ik} = \frac{1}{\sum_j \delta(Y^j = y_k)} \sum_j X_i^j \delta(Y^j = y_k)$$

ith feature

kth class

$\delta$(z)=1 if z true, else 0

$$\widehat{\sigma}_{ik}^2 = \frac{1}{\sum_j \delta(Y^j = y_k)} \sum_j (X_i^j - \widehat{\mu}_{ik})^2 \delta(Y^j = y_k)$$