

work with Data, instead of model

### 1.Learning → classification:

INPUT(X) — extraction →

(Eg.3D has x, y, z, color...)

Features—Machine learning→ predicted OUTPUT(Y)

### 2.ML:

train data(some domain):

attributes, class belongs to

(learning balanced way: learn and able to generate not matching)

tune parameter: not feature parameter, it optimizes ML model and training process

Eg.

Spam filter

SET UP: get a large collection, labeled each

Note: someone has to hand label all this data!

Want to learn to predict labels of new , future emails

INPUT: an email

FEATURE(attributes): spam/ham feature

OUTPUT: Spam/ham

Eg.

Digit Recognition

- Input: images / pixel grids
- Output: a digit 0-9

#### Setup:

- Get a large collection of example images, each labeled with a digit
- Note: someone has to hand label all this data
- Want to learn to predict labels of new, future digit images

- Features: the attributes used to make the digit decision
  - Pixels: (6,8)-ON
  - Shape Patterns: NumComponents, AspectRatio, NumLoops
  - ...



other classification also work with ML object recognition, Medical diagnosis..

## Model based Classification

### ■ Model-based approach

- Build a model (e.g. Bayes' net) where both the label and features are random variables
- Instantiate any observed features
- Query for the distribution of the label conditioned on the features

### ■ Challenges

- What structure should the BN have?
- How should we learn its parameters?

## Naïve Bayes Model

### ■ Random variables in this Bayes' net:

- The label

- $F_1, F_2, \dots, F_n$  = The n features

### ■ Probability tables in this Bayes' net:

- $P(Y)$  = Probability of each label occurring, given no information about the features. Sometimes called the prior.

- $P(F_i|Y)$  = One table per feature. Probability distribution over a feature, given the label.



### P(Y) baseline

## General Naïve Bayes

### ■ A general Naive Bayes model:

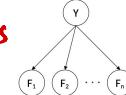
$\text{Y} \rightarrow \text{class}$   
 $F \rightarrow \text{Features}$

$$P(Y, F_1, \dots, F_n) = P(Y) \prod_i P(F_i|Y)$$

$|Y|$  parameters

$n \times |F| \times |Y|$

parameters



### ■ We only have to specify how each feature depends on the class

### ■ Total number of parameters is linear in n

### ■ Model is very simplistic, but often works anyway

## Example: Naïve Bayes for Spam Filter

### ■ Step 4.2: Inference

### ■ Instantiate features (evidence):

- $F_1 = \text{yes}$
- $F_2 = 1$

### ■ Compute joint probabilities:

- $P(Y = \text{spam} | F_1 = \text{yes}, F_2 = 1) = P(Y = \text{spam}) P(F_1 = \text{yes} | \text{spam}) P(F_2 = 1 | \text{spam})$
- $P(Y = \text{ham} | F_1 = \text{yes}, F_2 = 1) = P(Y = \text{ham}) P(F_1 = \text{yes} | \text{ham}) P(F_2 = 1 | \text{ham})$

### ■ Normalize:

- $0.4 \cdot 0.7 \cdot 0.7 = 0.294$
- $P(Y = \text{spam} | F_1 = \text{yes}, F_2 = 1) = 0.0048 / (0.0048 + 0.294) = 0.14$
- $P(Y = \text{ham} | F_1 = \text{yes}, F_2 = 1) = 0.294 / (0.0048 + 0.294) = 0.86$

### ■ Classification result:

- 14% chance the email is spam, 86% chance it's ham.
- Or, if you don't need probabilities, note that 0.0294 > 0.0048 and guess ham.

there maybe lots of words, Ignor Less Informative Words

## Inference for Naïve Bayes

### ■ Goal: compute posterior distribution over label variable Y

### ■ Step 1: get joint probability of label and evidence for each label

$$P(Y, f_1, \dots, f_n) = \begin{bmatrix} P(y_1, f_1, \dots, f_n) \\ P(y_2, f_1, \dots, f_n) \\ \vdots \\ P(y_k, f_1, \dots, f_n) \end{bmatrix} \xrightarrow{\text{vectors}} \begin{bmatrix} P(y_1) \prod_i P(f_i|y_1) \\ P(y_2) \prod_i P(f_i|y_2) \\ \vdots \\ P(y_k) \prod_i P(f_i|y_k) \end{bmatrix} \xrightarrow{+} P(f_1, \dots, f_n) \xrightarrow{+} P(Y|f_1, \dots, f_n)$$

### ■ Step 2: sum to get probability of evidence

### ■ Step 3: normalize by dividing Step 1 by Step 2

## Eg. NB

## Naïve Bayes for Digits

### ■ Naïve Bayes: Assume all features are independent effects of the label

### ■ Simple digit recognition version:

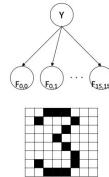
- One feature (variable)  $F_{i,j}$  for each grid position  $i,j$
- Feature values are on/off based on whether intensity is more or less than 0.5 in underlying image
- Each input maps to a feature vector, e.g.

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \rightarrow (F_{0,0} = 0, F_{0,1} = 0, F_{0,2} = 1, F_{0,3} = 1, F_{0,4} = 0, \dots, F_{15,15} = 0)$$

Here: lots of features, each is binary valued

$$\text{■ Naïve Bayes model: } P(Y|F_{0,0} \dots F_{15,15}) \propto P(Y) \prod_{i,j} P(F_{i,j}|Y)$$

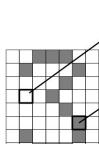
### ■ What do we need to learn?



## Naïve Bayes for Digits: Parameters

### P(Y)

1	0.1
2	0.1
3	0.1
4	0.1
5	0.1
6	0.1
7	0.1
8	0.1
9	0.1
0	0.1



1	0.01
2	0.05
3	0.05
4	0.30
5	0.80
6	0.90
7	0.05
8	0.60
9	0.50
0	0.80

1	0.05
2	0.01
3	0.90
4	0.80
5	0.90
6	0.90
7	0.25
8	0.85
9	0.60
0	0.80

## General Naïve Bayes

$$Posterior \text{ Likelihood} = \frac{P(E|H) \cdot P(H)}{P(E)} \cdot \frac{\text{Prior}}{\text{Prior}}$$

这张幻灯片介绍了使用朴素贝叶斯分类器所需的组件以及这些组件的来源。

- 推理方法:** 这部分我们刚刚讨论过。使用朴素贝叶斯分类器进行推理时，我们首先需要一组概率值：类别的先验概率  $P(Y)$  和条件概率  $P(F_i|Y)$ 。然后，我们使用标准的推理方法来计算给定所有特征  $F_1, F_2, \dots, F_n$  后，类别  $Y$  的后验概率  $P(Y|F_1, \dots, F_n)$ 。这里并没有什么新颖的内容，这是朴素贝叶斯方法的标准操作。

- 局部类别概率表的估计:** 我们需要两种类型的概率估计值：

- 类别的先验概率  $P(Y)$ ，也就是在观察任何特征之前，每个类别可能出现的概率。  
important  
这些概率值被统称为模型的参数（这里用符号  $\theta$  表示）。
- 每个特征给定类别的条件概率  $P(F_i|Y)$ ，也就是当我们知道某个类别时，每个特征出现的概率。

在之前的讨论中，我们可能会认为这些概率值凭空出现，但实际上它们通常通过对训练数据进行统计而获得的。例如，通过计算训练数据集中每个类别出现的频率来估计先验概率，以及在每个类别下，不同特征出现的频率来估计条件概率。

## ERM: no classifier

这张幻灯片主要讨论了如何评估机器学习模型的质量，以及“经验风险最小化”(Empirical Risk Minimization, ERM) 的概念。

- 模型质量评估:** 在机器学习中，我们通常想要评估模型的质量，确保它在未知数据上表现良好。

理想情况下，我们希望模型在真实的测试数据分布上表现良好。

- 经验风险最小化:** ERM是机器学习的一个基本原则。它的目标是找到在训练数据集上表现最佳的模型。由于我们通常不知道数据的真实分布，ERM策略就是选择在实际训练集上“看起来”最好的模型。将寻找最佳模型的过程描述为一个优化问题。目标是最小化训练集上的损失函数，这个损失函数衡量的是模型预测和实际值之间的误差。

- 主要挑战:** 拟合

- 过拟合是指模型在训练数据上表现非常好，但在未见过的数据上表现不佳。这通常发生在模型过于复杂，以至于它学会了训练数据中的噪声，而不仅仅是底层的数据生成分布。
- 为了缓解过拟合，增加训练数据是一个有效的策略。因为更多的数据可以减少采样的方差，并使训练数据更好地代表真实的测试数据分布。
- 另一个策略是限制模型假设的复杂度，比如通过正则化技术，这可以鼓励模型简化它的结构，减少对训练数据中噪声的学习。

## Leading to Parameter Estimate (Distribution of Variable)

Elicitation: ask human

Empirical way: use training data

Eg.

- Example: The parameter  $\theta$  is the true fraction of red beans in the jar. You don't know  $\theta$  but would like to estimate it.
- Collecting training data: You randomly pull out 3 beans:
- Estimating  $\theta$  using counts, you guess 2/3 of beans in the jar are red.
- Can we mathematically show that using counts is the "right" way to estimate  $\theta$ ?

right answer?

proof:

### Parameter Estimation with Maximum Likelihood

- Maximum likelihood estimation:** Choose the  $\theta$  value that maximizes the probability of observation

In other words, choose the  $\theta$  value that maximizes  $P(\text{observation} | \theta)$

For our problem:

(red)

Randomly selected 2 red and 1 blue | 0 beans are red)

→ P(red |  $\theta$ ) \* P(blue |  $\theta$ )

→  $P(1|0)$

We want to compute:

$\arg\max_{\theta} \theta^0 (1-\theta)^1$

We want to compute:

$\arg\max_{\theta} \theta^0 (1-\theta)^1$  If need, math way

Set derivative to 0, and solve!

Common mistake: the likelihood (expression we're maxing) is the product of a lot of probabilities. This can make derivatives hard to work with.

Solution: Maximize the log-likelihood instead. Useful fact:

$\arg\max_{\theta} f(\theta) = \arg\max_{\theta} \ln f(\theta)$

## Important Concepts

- How do we check that we're not overfitting during training?

Split training data into 3 different sets:

- Training set
- Held out set (more on this later)
- Test set

- Experimentation cycle

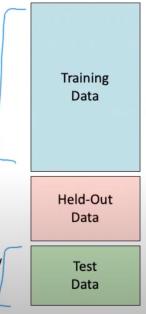
- Learn parameters (e.g. model probabilities) on training set
- Compute accuracy of test set
- Very important: never “peek” at the test set!

- Evaluation (many metrics possible, e.g. accuracy)

- Accuracy: fraction of instances predicted correctly

- Overfitting and generalization

- Want a classifier which does well on test data
- Overfitting: fitting the training data very closely, but not generalizing well
- We'll investigate overfitting and generalization formally in a few lectures



We don't want to be too specific, but smooth and regularize the estimate

## Parameter Estimation

- Estimating the distribution of a random variable

- Elicitation: ask a human (why is this hard?)

- Empirically: use training data (learning)

- E.g.: for each outcome  $x$ , look at the empirical rate of that value:

$$P_{ML}(x) = \frac{\text{count}(x)}{\text{total samples}} \rightarrow \text{Count } r \quad r \quad r \quad b \quad \rightarrow 2r+1b \quad P_{ML}(r) = 2/3$$



- This is the estimate that maximizes the likelihood of the data

$$L(x, \theta) = \prod_i P_{\theta}(x_i)$$

## Maximum Likelihood?

- Relative frequencies are the maximum likelihood estimates

$$\theta_{ML} = \arg \max_{\theta} P(X|\theta) \rightarrow \text{features} \quad \text{not ML itself}$$

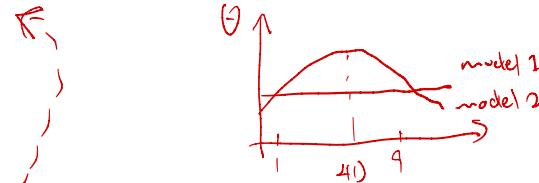
$$= \arg \max_{\theta} \prod_i P_{\theta}(X_i) \quad P_{ML}(x) = \frac{\text{count}(x)}{\text{total samples}}$$

- Another way showing it (semi proof) is as follows:

$$\arg \max_{\theta} P(\theta|X) = \arg \max_{\theta} \prod_i P_{\theta}(X_i) \rightarrow \text{by Bayes rule}$$

$$= \arg \max_{\theta} P(X|\theta) P(\theta) / P(X) \rightarrow ???$$

$$= \arg \max_{\theta} P(X|\theta) P(\theta)$$



# Parameter Estimation with Maximum Likelihood

$\underset{\theta}{\operatorname{argmax}} \theta^2(1-\theta)$	Find $\theta$ that maximizes likelihood
$= \underset{\theta}{\operatorname{argmax}} \ln(\theta^2(1-\theta))$	Find $\theta$ that maximizes log-likelihood (will be the same $\theta$ )
$\frac{d}{d\theta} \ln(\theta^2(1-\theta)) = 0$	Set derivative to 0
$\frac{d}{d\theta} [\ln(\theta^2) + \ln(1-\theta)] = 0$	Logarithm rule: products become sums
$\frac{d}{d\theta} [2\ln(\theta) + \ln(1-\theta)] = 0$	Logarithm rule: exponentiation becomes multiplication
$\frac{d}{d\theta} 2\ln(\theta) + \frac{d}{d\theta} \ln(1-\theta) = 0$	Now we can derive each term of the original product
$\frac{2}{\theta} - \frac{1}{1-\theta} = 0$	separately Reminder: Derivative of $\ln(\theta)$ is $1/\theta$
$\theta = \frac{2}{3}$	Use algebra to solve for $\theta$ . If we used arbitrary red and blue counts $r$ and $b$ instead of $r=2$ and $b=1$ , we'd get $\theta = r / (r+b)$ , the count estimate.

## Linear Vector

$x \rightarrow f(x) \rightarrow y$   
letter  $\rightarrow$  [features]  $\rightarrow$  Spam or Ham

## Laplace Smoothing

- Laplace's estimate (extended):
  - Pretend you saw every outcome  $k$  extra times



$$P_{LAP,k}(x) = \frac{c(x) + k}{N + k|X|}$$

$$P_{LAP,0}(X) =$$

- What's Laplace with  $k = 0$ ?
- $k$  is the strength of the prior

$$P_{LAP,1}(X) =$$

$$P_{LAP,100}(X) =$$

- Laplace for conditionals:

$$\begin{aligned} P_{LAP,k}(x|y) &= \frac{c(x,y) + k}{c(y) + k|X|} \\ &\text{atly:} \end{aligned}$$

- $c(x)$  是观测到的事件  $x$  的次数。
- $N$  是总的观测次数。
- $|X|$  是所有不同事件的数量。
- $k$  是平滑参数，也称作先验的强度。

在拉普拉斯平滑中，如果  $k = 0$ ，则没有平滑效果。我们得到的是事件的经验概率（直接的频率概率）。当  $k > 0$  时， $k$  越大，平滑效果越明显，也就是说对概率估计的影响越大。例如， $k = 1$  时通常被称为加一平滑（因为给每个事件加上 1）。这在许多情况下可以避免零概率问题，并且提供了一个比较温和的平滑效果。而  $k = 100$  会给予先验一个非常强的权重，可能会过度平滑，使得估计概率偏离实际观察的数据。

## Error

- Need more features – words aren't enough!
  - Have you emailed the sender before?
  - Have 1K other people just gotten the same email?
  - Is the sending information consistent?
  - Is the email in ALL CAPS?
  - Do inline URLs point where they say they point?
  - Does the email address you by (your) name?
- Can add these information sources as new variables in the NB model