# Machine Learning

- Machine learning: how to acquire a model from data / experience
  - Learning parameters (e.g. probabilities)
  - Learning structure (e.g. BN graphs)
  - Learning hidden concepts (e.g. clustering, neural nets)
- Today: model-based classification with Naive Bayes

#### Classification



### **General Naïve Bayes**

• A general Naive Bayes model:

|Y| parameters

$$P(\mathbf{Y}, \mathbf{F}_1 \dots \mathbf{F}_n) = P(\mathbf{Y}) \prod_i P(\mathbf{F}_i | \mathbf{Y})$$

$$|\mathbf{Y}| \ge |\mathbf{F}|^n \text{ values}$$

$$n \ge |\mathbf{F}| \ge |\mathbf{Y}|$$

n x |F| x |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



### **General Naïve Bayes**

- What do we need in order to use Naïve Bayes?
  - Inference method (we just saw this part)
    - Start with a bunch of probabilities: P(Y) and the P(F<sub>i</sub>|Y) tables
    - Use standard inference to compute P(Y|F<sub>1</sub>...F<sub>n</sub>)
    - Nothing new here
  - Estimates of local conditional probability tables
    - P(Y), the prior over labels
    - P(F<sub>i</sub>|Y) for each feature (evidence variable)
    - These probabilities are collectively called the *parameters* of the model and denoted by  $\theta$
    - Up until now, we assumed these appeared by magic, but...
    - ...they typically come from training data counts: we'll look at this soon

# Naïve Bayes for Digits

- Naïve Bayes: Assume all features are independent effects of the label
- Simple digit recognition version:
  - One feature (variable) F<sub>ii</sub> 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.

$$\rightarrow \langle F_{0,0} = 0 \ F_{0,1} = 0 \ F_{0,2} = 1 \ F_{0,3} = 1 \ F_{0,4} = 0 \ \dots F_{15,15} = 0 \rangle$$

- Here: lots of features, each is binary valued
- Naïve Bayes model:  $P(Y|F_{0,0} \dots F_{15,15}) \propto P(Y) \prod P(F_{i,j}|Y)$  (Bayes' theorem)
- What do we need to learn?



# Naïve Bayes for Text

- Bag-of-words Naïve Bayes:
  - Features: W<sub>i</sub> is the word at position i
  - As before: predict label conditioned on feature variables (spam vs. ham)
  - As before: assume features are conditionally independent given label
  - New: each W<sub>i</sub> is identically distributed

Word at position *i*, not i<sup>th</sup> word in the dictionary!

• Generative model: 
$$P(Y, W_1 \dots W_n) = P(Y) \prod_i P(W_i | Y)$$

- "Tied" distributions and bag-of-words
  - Usually, each variable gets its own conditional probability distribution P(F|Y)
  - In a bag-of-words model
    - Each position is identically distributed
    - All positions share the same conditional probs P(W|Y)
    - Why make this assumption?
  - Called "bag-of-words" because model is insensitive to word order or reordering

# Important Concepts

- Data: labeled instances (e.g. emails marked spam/ham)
  - Training set
  - Held out set
  - Test set
- Features: attribute-value pairs which characterize each x
- Experimentation cycle
  - Learn parameters (e.g. model probabilities) on training set
  - (Tune hyperparameters on held-out 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



# **Empirical Risk Minimization**

#### Empirical risk minimization

- Basic principle of machine learning
- We want the model (classifier, etc) that does best on the true test distribution
- Don't know the true distribution so pick the best model on our actual training set
- Finding "the best" model on the training set is phrased as an optimization problem

#### Main worry: overfitting to the training set

- Better with more training data (less sampling variance, training more like test)
- Better if we limit the complexity of our hypotheses (regularization and/or small hypothesis spaces)

### **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:

$$P_{\mathsf{ML}}(x) = \frac{\mathsf{count}(x)}{\mathsf{total samples}}$$

$$\begin{array}{c} \mathbf{r} \quad \mathbf{r} \quad \mathbf{b} \\ P_{\mathsf{ML}}(\mathbf{r}) = 2/3 \end{array}$$



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

$$\theta_{ML} = \arg \max_{\theta} P(\mathbf{X}|\theta)$$
  
=  $\arg \max_{\theta} \prod_{i} P_{\theta}(X_{i})$   $P_{\mathsf{ML}}(x) = \frac{\operatorname{count}(x)}{\operatorname{total samples}}$ 

Another option is to consider the most likely parameter value given the data

$$\theta_{MAP} = \arg \max_{\theta} P(\theta | \mathbf{X})$$
  
=  $\arg \max_{\theta} P(\mathbf{X} | \theta) P(\theta) / P(\mathbf{X})$   $\longrightarrow$  ????  
=  $\arg \max_{\theta} P(\mathbf{X} | \theta) P(\theta)$ 

# **Example: Overfitting**



2 wins!!

# Overfitting



#### **Unseen Events**



# Laplace Smoothing

- Laplace's estimate:
  - Pretend you saw every outcome once more than you actually did

$$P_{LAP}(x) = \frac{c(x) + 1}{\sum_{x} [c(x) + 1]}$$
$$= \frac{c(x) + 1}{N + |X|}$$

$$P_{ML}(X) =$$

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

 Can derive this estimate with Dirichlet priors (see cs281a)

# Laplace Smoothing

#### Laplace's estimate (extended):

Pretend you saw every outcome k extra times

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

- What's Laplace with k = 0?
- k is the strength of the prior
- Laplace for conditionals:
  - Smooth each condition independently:

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

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

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

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

### Estimation: Linear Interpolation\*

- In practice, Laplace often performs poorly for P(X|Y):
  - When |X| is very large
  - When |Y| is very large
- Another option: linear interpolation
  - Also get the empirical P(X) from the data
  - Make sure the estimate of P(X|Y) isn't too different from the empirical P(X)

$$P_{LIN}(x|y) = \alpha \hat{P}(x|y) + (1.0 - \alpha)\hat{P}(x)$$

- What if α is 0? 1?
- This is called 'backoff'

# **Real NB: Smoothing**

- For real classification problems, smoothing is critical
- New odds ratios:

 $rac{P(W| extsf{spam})}{P(W| extsf{ham})}$ 

helvetica	:	11.4
seems	:	10.8
group	:	10.2
ago	:	8.4
areas	:	8.3

P(W|ham)

 $\overline{P(W|\text{spam})}$ 

verdana	:	28.8
Credit	:	28.4
ORDER	•	27.2
<font></font>	:	26.9
money	:	26.5
•••		



Do these make more sense?

# Tuning



# Tuning on Held-Out Data

#### Now we've got two kinds of unknowns

- Parameters: the probabilities P(X|Y), P(Y)
- Hyperparameters: e.g. the amount / type of smoothing to do, k, α

#### What should we learn where?

- Learn parameters from training data
- Tune hyperparameters on different data
  - Why?
- For each value of the hyperparameters, train and test on the held-out data
- Choose the best value and do a final test on the test data



#### Features



#### Errors, and What to Do

#### Examples of errors

Dear GlobalSCAPE Customer,

GlobalSCAPE has partnered with ScanSoft to offer you the latest version of OmniPage Pro, for just \$99.99\* - the regular list price is \$499! The most common question we've received about this offer is - Is this genuine? We would like to assure you that this offer is authorized by ScanSoft, is genuine and valid. You can get the . . .

. . . To receive your \$30 Amazon.com promotional certificate, click through to

http://www.amazon.com/apparel

and see the prominent link for the \$30 offer. All details are there. We hope you enjoyed receiving this message. However, if you'd rather not receive future e-mails announcing new store launches, please click . . .

# What to Do About Errors?

#### 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
- Next class we'll talk about classifiers which let you easily add arbitrary features more easily, and, later, how to induce new features



#### Baselines

#### • First step: get a baseline

- Baselines are very simple "straw man" procedures
- Help determine how hard the task is
- Help know what a "good" accuracy is
- Weak baseline: most frequent label classifier
  - Gives all test instances whatever label was most common in the training set
  - E.g. for spam filtering, might label everything as ham
  - Accuracy might be very high if the problem is skewed
  - E.g. calling everything "ham" gets 66%, so a classifier that gets 70% isn't very good...
- For real research, usually use previous work as a (strong) baseline

## **Confidences from a Classifier**

#### • The confidence of a probabilistic classifier:

Posterior probability of the top label

 $confidence(x) = \max_{y} P(y|x)$ 

- Represents how sure the classifier is of the classification
- Any probabilistic model will have confidences
- No guarantee confidence is correct

#### Calibration

- Weak calibration: higher confidences mean higher accuracy
- Strong calibration: confidence predicts accuracy rate
- What's the value of calibration?



### Summary

- Bayes rule lets us do diagnostic queries with causal probabilities
- The naïve Bayes assumption takes all features to be independent given the class label
- We can build classifiers out of a naïve Bayes model using training data
- Smoothing estimates is important in real systems
- Classifier confidences are useful, when you can get them