Introduction to Machine Learning

roduction to Machine Learning Amo G. Tons

Amo G. Tong

Lecture 5 Supervised Learning

- · Naïve Bayes Classifier
- · Readings: Mitchell Ch 6.9-6.10; Murphy Ch 3.5
- · Some materials are courtesy of Vibhave Gogate and Tom Mitchell.
- All pictures belong to their creators.

Introduction to Machine Learning

mo G. Tong

Step 3. Select a measure to evaluate the functions in H. Step 4. Use a machine learning algorithm to find the be-

• Step 1. Select the features of x to be used.

 Step 4. Use a machine learning algorithm to find the best function in H according to your measure.

Step 2. Select a hypothesis space H: a set of candidate functions to

• Given some training examples < x, f(x) > and an unknown function f.

- · Bayesian Learning: the most probable classification
- · Naïve Bayes Classifier: additional assumption.

duction to Machine Learning

approximate f.

Supervised Learning

• Find a good approximation of f.

Amo G. Tong

Optimal Bayes Classifier

- Instance space X: each $x \in X$ is characterized by attributes $(x_1, ..., x_k)$
- Classifications V.

Sky	Water	Forecast	EnjoySport
Sunny	Warm	Same	Yes

• Underlying truth: a distribution Pr[x, v] over (x, v).

Most probable classification = $\operatorname{argmax}_{v \in V} \Pr[v|x]$

Introduction to Machine Learning

Amo G. Tong

Optimal Bayes Classifier

- Settings: learn a function $X \to V$
- Instance space X: each $x \in X$ is characterized by attributes $(x_1, ..., x_k)$
- Classifications V.

Sky	Water	Forecast	EnjoySport
Sunny	Warm	Same	Yes

$$\operatorname{argmax}_{v \in V} \Pr[v|x] = \operatorname{argmax}_{v \in V} \sum_{h \in H \text{ and } h(x) = v} \Pr[h]$$

The probability that x is classified as v.

x=(Sunny, Warm, Same)

Pr[yes|x] is the probability that x is classified as yes.

Introduction to Machine Learnin

Amo G. Ton

Optimal Bayes Classifier

- Settings: learn a function $X \to V$
- Instance space X: each $x \in X$ is characterized by attributes $(x_1, ..., x_k)$
- Classifications V.

Sky	Water	Forecast	EnjoySport
Sunny	Warm	Same	Yes

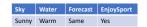
 $\operatorname{argmax}_{v \in V} \Pr[v|x] = \operatorname{argmax}_{v \in V} \sum_{h \in H \text{ and } h(x) = v} 1$

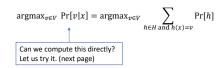
The distribution over H. This is hard to compute. We need to estimate $\Pr[h]$ by $\Pr[h|D]$ for each h in H. The example in Lec 4 is very simple.

Introduction to Machine Learning Amo G. Tong

Optimal Bayes Classifier

- Settings: learn a function $X \to V$
- Instance space X: each $x \in X$ is characterized by attributes $(x_1, ..., x_k)$
- Classifications V.





Amo G. Tong

Naïve Bayes Classifier

$$Pr[v|x] = \frac{Pr[x|v] Pr[v]}{Pr[x]}$$
 Bayes theorem.

 $\operatorname{argmax} \Pr[v|x] = \operatorname{argmax} \ \Pr[x|v] \Pr[v] = \operatorname{argmax} \ \Pr[x_1, x_2, \dots, x_k|v] \Pr[v]$

$$x = (x_1, \dots, x_k)$$



Can you classify a new instance (Sunny)? $\Pr[\operatorname{Sunny}|yes,D]\Pr[yes|D] = \frac{1}{2} * \frac{2}{5}$ $\Pr[\operatorname{Sunny}|no,D]\Pr[no|D] = \frac{2}{2} * \frac{3}{5}$

to Machine Learning Amo G. Tong 8

Naïve Bayes Classifier

$$Pr[v|x] = \frac{Pr[x|v] Pr[v]}{Pr[x]}$$
 Bayes theorem.

 $\operatorname{argmax} \Pr[v|x] = \operatorname{argmax} \ \Pr[x|v] \ \Pr[v] = \operatorname{argmax} \ \Pr[x_1, x_2, \dots, x_k|v] \ \Pr[v]$ $x = (x_1, \dots, x_k)$

Suppose x_i and y are binary.

There are totally $\Theta(2^k)$ different parameters (instance, classification).

A lot of data needed.

Introduction to Machine Learning Amo G. Tong

Naïve Bayes Classifier

$$\Pr[v|x] = \frac{\Pr[x|v] \Pr[v]}{\Pr[x]} \quad \text{Bayes theorem}.$$

 $\operatorname{argmax} \Pr[v|x] = \operatorname{argmax} \Pr[x|v] \Pr[v] = \operatorname{argmax} \Pr[x_1, x_2, \dots, x_k|v] \Pr[v]$ $x = (x_1, \dots, x_k)$

Suppose x_i and y are binary.

There are totally $\Theta(2^k)$ different parameters (instance, classification).

Sky	Water	Forecast	EnjoySport
Sunny	Warm	Same	Yes
Sunny	Warm	Same	No
Rainy	Warm	Change	No
Sunny	Cool	Change	Yes
Sunny	Cool	Same	no

A lot of data needed.

Can you classify a new instance (Sunny, Warm, Change)? Pr[Sunny, Warm, Change|yes, D]=0 Pr[Sunny, Warm, Change | no, D] = 0

Introduction to Machine Learning

Amo G. Tong

Naïve Bayes Classifier

 $\operatorname{argmax} \Pr[v|x] = \operatorname{argmax} \ \Pr[x|v] \Pr[v] = \operatorname{argmax} \ \Pr[x_1, x_2, \dots, x_k|v] \Pr[v]$ $x = (x_1, \dots, x_k)$

Suppose x_i and y are binary.

There are totally $\Theta(2^k)$ different parameters (instance, classification).

Naïve Bayes Assumption: given the classification v, the attributes are independent.

 $Pr[x_1, x_2|v] = Pr[x_1|v]Pr[x_2|v]$

Naïve Bayes Classifier

Naïve Bayes Assumption: given the classification v, the attributes are independent.

 $Pr[x_1, x_2|v] = Pr[x_1|v]Pr[x_2|v]$

Sky Water Forecast EnjoySport Sunny Warm Same

Can you classify a new instance (Sunny, Warm, Change)? Pr[Sunny, Warm, Change|yes]

= Pr[Sunny|yes] * Pr[Warm|yes]* Pr[Change|yes]

Amo G. Tong

Naïve Bayes Classifier

 $\operatorname{argmax} \Pr[v|x] = \operatorname{argmax} \Pr[x|v] \Pr[v]$ = $\operatorname{argmax} \Pr[x_1, x_2, ..., x_k | v] \Pr[v]$

Naïve Baves Assumption: given the classification v, the attributes are independent.

 $\Pr[x_1,x_2|v] = \Pr[x_1|v] \Pr[x_2|v]$

- $\operatorname{argmax} \Pr[x_1, x_2, ..., x_k | v] \Pr[v] = \operatorname{argmax} \prod \Pr[x_i | v] \Pr[v]$
- Naïve Bayes Classifier: $\operatorname{argmax} \prod \Pr[x_i | v] \Pr[v]$

How many parameters now? O(k)

Introduction to Machine Learning

Amo G. Tong

Naïve Bayes Classifier

- $\operatorname{argmax} \Pr[x_1, x_2, ..., x_k | v] \Pr[v] = \operatorname{argmax} \prod \Pr[x_i | v] \Pr[v]$
- Naïve Bayes Classifier: $\operatorname{argmax} \prod \Pr[x_i | v] \Pr[v]$

Can you classify a new instance (Sunny, Warm, Change)?

Pr[Sunny|yes] * Pr[Warm|yes]* Pr[Change|yes]*Pr[yes]

Pr[Sunny|no] * Pr[Warm|no] * Pr[Change|no] * Pr[no]

Introduction to Machine Learning

Amo G. Tong

Naïve Bayes Classifier

- Naïve Bayes Classifier: $\operatorname{argmax} \prod \Pr[x_i \mid v] \Pr[v]$
- For a new instance $(x_1, ..., x_k)$ and each $v \in V$.
- Estimate $Pr[x_i | v]$ and Pr[v] for each x_i and v by data D.

MLE of Pr[v] and $Pr[x_i|v]$ Bernoulli Point Estimation.

- $C_{class=v}$: set of instances with classification v in D.
- $C_{X_i=x_i}$: set of instances with x_i as the *i*-th attribute.

 $\Pr[v]_{MLE} = \frac{|c_{class=v}|}{|D|} \qquad \Pr[x_i|v]_{MLE} = \frac{|C_{class=v} \cap C_{X_i=x_i}|}{|C_{class=v}|}$

Introduction to Machine Learning

Amo G. Tong

Naïve Bayes Classifier

- Naïve Bayes Classifier: $\operatorname{argmax} \prod \Pr[x_i \mid v] \Pr[v]$
- $\Pr[v] = \frac{|c_{class=v}|}{|D|}$, $\Pr[x_i|v] = \frac{|c_{class=v} \cap c_{X_i=x_i}|}{|c_{class=v}|}$

· How to classify (Sunny, Warm, Change)?

Pr[ves]=2/5, Pr[no] = 3/5,Pr[Sunny|yes] = 2/2,Pr[Sunny|no] = 2/3,Pr[Warm|yes] = 1/2Pr[Warm|no] = 2/3Pr[Change|yes] = 1/2Pr[Change|no] = 1/3Introduction to Machine Learning Amo G. Tong

Text Classification: Rumor or not Rumor

- · Article: a sequence of words.
- $X \rightarrow V$

- x = 'I have a red apple and a blue apple'
- X: all possible articles
- *V*={rumor, not rumor}
- D: some labeled articles

 $v_{opt} = \operatorname{argmax} \Pr[x|v] \Pr[v]$

Text Classification: Rumor or not Rumor

- Text Classification. Rumor or not Rumor??
- Article x: a sequence of words.

x= 'I have a red apple and a blue apple'

• $v_{ont} = \operatorname{argmax} \Pr[x|v] \Pr[v]$

There are so many sequences of words.

- · An article usually has more than 1000 words.
- · More than 10.000 common words.
- That is $10,000^{1000} = 10^{4000}$
- Atoms in universe: 10⁸⁰
- Hm
- We need some assumption/simplification.

Introduction to Machine Learning

Amo G. Tong

Text Classification: Rumor or not Rumor

- Text Classification. Rumor or not Rumor??
- Article x: a sequence of words. (L: length)

 $\operatorname{argmax}_{v} \Pr[x|v] \Pr[v]$

x= 'I have a red apple and a blue apple'

• Naïve Bayes Assumption: the positions are independent given the class.

$$\operatorname{argmax}_{v} \prod_{x=1}^{L} \Pr[X_i = x_i | v] \Pr[v]$$

• $Pr[X_i = x_i | v]$ the probability of observing x_i in the *i*-th position.

I | have | a | red | apple | and | a | blue | apple

Introduction to Machine Learning

Amo G. Tong

Text Classification: Rumor or not Rumor

- $Pr[X_i = x_i | v]$ the probability of observing x_i in the *i*-th position.
- How to compute $Pr[X_i = x_i | v]$?
- · Positions have the same distribution over the word.
- $Pr[X_i = apple | v] = Pr[X_i = apple | v]$
- $Pr[X_i = apple | v]$ the probability that the *i*-th position is 'apple'.
- $\Pr[x_i|v]$ can be computed by $\frac{Count(x_i,v)}{\sum_{x\in X}Count(x,v)}$ where X is the set of all distinct words, and count(x,v) is number of the times word x appears in documents of class v.

Remarks

- What if part of the new instance is never observed?
- Suppose you have $x = (x_1, cool, x_3)$ but no training instance has Water=Cool.
- We will have $\prod \Pr[x_i | yes] \Pr[yes] = \prod \Pr[x_i | no] \Pr[no] = 0$ regardless of other attributes of x.

 $\Pr[x_i|v]_{MLE} = \frac{\left|C_{Y=v} \cap C_{X_i=x_i}\right|}{\left|C_{Y=v}\right|}$

· This is not reasonable.

Amo G. Tong

Remarks

- What if part of the new instance is never observed?
- · Let us use 'smoothing'.

 $\Pr[x_i|v]_{MLE} = \frac{\left|C_{class=v} \cap C_{X_i=x_i}\right| + mp}{\left|C_{class=v}\right| + m}$

- · m-estimate
- m is a constant and p is the prior estimate of x_i over possible values. If no prior estimate, set p = 1/k where k is the number of possible values of the attribute.
- Pretend you have m extra instances of class v and mp of them have x_i as the attribute value.

Introduction to Machine Learning

Amo G. Tong

Remarks

- What if part of the new instance is never observed?
- Suppose you have $x = (x_1, cool, x_3)$ but no training instance has Water=Cool.

$$\Pr[x_i|v]_{MLE} = \frac{\left|C_{class=v} \cap C_{X_i=x_i}\right| + mp}{\left|C_{class=v}\right| + m}$$

{cool, warm}

$Pr[cool yes] = \frac{ C_{Y=v} + C_{X_i=x_i} + mp}{ C_{Y=v} + m} = \frac{m + 0.3}{2 + m}$	Sunny
	Sunny
$Pr[cool no] = \frac{ C_{Y=v} \cap C_{X_i=x_i} + mp}{ C_{v-1} + m} = \frac{m * 0.5}{2 + m}$	Rainy
$ C_{Y=v} + m \qquad 2 + m$	Sunny

Amo G. Tong

Example

A text classification example (from Dan Jurafsky)

	Doc	Words	Class
Training	1	Chinese Beijing Chinese	С
	2	Chinese Chinese Shanghai	С
	3	Chinese Macao	С
	4	Tokyo Japan Chinese	j
Test	5	Chinese Chinese Tokyo Japan	?

- We use m-estimate where (m=number of distinct word) and (p=1/m).
- m=6, p=1/6

Pr[j]=1/4 Pr[c]=3/4 Pr[Chinese | i]=(1+1)/(3+6)=2/9 Pr[Tokyo|j]=(1+1)/(3+6)=2/9 Pr[Japan|j]=(1+1)/(3+6)=2/9

• $\Pr[j|x] = \left(\frac{1}{4}\right) * \left(\frac{2}{9}\right)^3 * \left(\frac{2}{9}\right) * \left(\frac{2}{9}\right) \approx 0.0001$ • $\Pr[c|x] = \left(\frac{3}{4}\right) * \left(\frac{3}{14}\right) * \left(\frac{1}{14}\right) * \left(\frac{1}{14}\right) \approx 0.0003$

Pr[Chinese | c]=(5+1)/(8+6)=3/7 Pr[Tokyo|c]=(0+1)/(8+6)=1/14 Pr[Japan|c]=(0+1)/(8+6)=1/14

Introduction to Machine Learning

Amo G. Tong

Remarks

- Naïve Bayes Assumption: the attributes are conditionally independent.
- · What if the assumption is not true?

$$\Pr[x_1, x_2, \dots, x_k | v] \Pr[v] \neq \prod \Pr[x_i | v]$$

- This assumption may not be true, but this approximate performs good.
- A plausible reason: Only need the probability of the correct class to be the largest!
- E.g.

Introduction to Machine Learning

- Truth: Pr[f(x) = 1] = 0.99 and Pr[f(x) = 0] = 0.01. Very likely f(x) = 1
- Your estimate: Pr[f(x) = 1] = 0.51 and Pr[f(x) = 0] = 0.49. You say f(x) = 1.

Remarks

- Naïve Bayes Assumption: the attributes are conditionally independent.
- · The effect of this assumption:
- (a) the hypothesis space is restricted to those satisfy this assumption.
- (b) it requires less data (low variance).

$$\Pr[x_1, x_2, \dots, x_k | v] \qquad \Pr[x_i | v]$$

(Sky=Sunny, Temperature=Warm, Normal, Strong, Warm, Same)=Yes

- (c) Naïve Bayes classifier finds the most probable classification within this hypothesis.
- If the assumption is true, Naïve Bayes is the optimal classifier.

Dealing with Small Numbers

- A practical issue: $\prod \Pr[x_i | v]$
- · We are multiplying lots of small numbers: underflow.
- $0.5^{57} = 7E 18$
- $\Pr[j|x] = \left(\frac{1}{4}\right) * \left(\frac{2}{9}\right)^3 * \left(\frac{2}{9}\right) * \left(\frac{2}{9}\right) \approx 0.0001$
- $\Pr[c|x] = \left(\frac{3}{4}\right) * \left(\frac{3}{7}\right)^3 * \left(\frac{1}{14}\right) * \left(\frac{1}{14}\right) \approx 0.0003$
- · Solution: use log and add.
- $a * b = e^{\ln a + \ln b}$
- Keep the log form.

Introduction to Machine Learning

Amo G. Tong

Gaussian Naïve Bayes

· Continuous features

Sky	Water	Forecast	EnjoySport
Sunny	Warm	Same	Yes

Sky	Water	Forecast	EnjoySport
0.2	0.5	0.7	Yes

$$\operatorname{argmax}_{\mathbf{v}} \prod_{x=1}^{k} \Pr[X_i = x_i | v] \Pr[v]$$

• $Pr[v]_{MLE}$ can be computed in the same way.

•
$$Pr[Sky = sunny|yes]_{MLE} = \frac{yes \wedge sunny}{yes}$$

• $Pr[Sky = 0.2|yes]_{MLE} = ??$

Generally, assume a distribution over the domain of Sky.

- · Discrete domain {sunny, cloudy} · Continuous domain R.
- Assume that Pr[Sky = x|yes] follows a particular distribution.

Introduction to Machine Learning

Amo G. Tong

Gaussian Naïve Bayes

· Continuous features

Sky	Water	Forecast	EnjoySport
0.2	0.5	0.7	Yes

- Assume Pr[Sky = x|yes] follows a particular distribution.
- Naïve Bayes + Gaussian
- X_i : the *i*-th attribute, v_k the k-th classification.
- Assume $Pr[X_i = x_i | v_k]$ follows $N(\mu_{i,k}, \sigma_{i,k}^2)$

$$\Pr[X_i = x | v_k] = \frac{1}{\sigma_{i,k} \sqrt{2\pi}} e^{\frac{-(x_i - \mu_{i,k})^2}{2\sigma_{i,k}^2}}$$

Point estimation:

Introduction to Machine Learning

 $\Pr[X_i = x | v_k] = \frac{1}{\sigma_{i,k} \sqrt{2\pi}} e^{\frac{-(x_i - \mu_{i,k})^2}{2\sigma_{i,k}^2}}$

Gaussian Naïve Bayes

- · Continuous features
- X_i : the *i*-th attribute, v_k the *k*-th class.
- Assume $\Pr[X_i = x_i | v_k]$ follows $N(\mu_{i,k}, \sigma_{i,k}^2)$

Sometimes we assume: The variance is Independent of class $\sigma_{i,k}^2 = \sigma_i^2$ Independent of attribute $\sigma_{i,k}^2 = \sigma_k^2$ Independent of both $\sigma_{ik}^2 = \sigma^2$

Point estimation:
Learn
$$\mu_{i,k}$$
, $\sigma_{i,k}^2$ from data

Amo G. Tong

Gaussian Naïve Bayes

Continuous features

$$\Pr[X_i = x_i | v_k] = \frac{1}{\sigma_{i,k} \sqrt{2\pi}} e^{\frac{-(x_i - \mu_{i,k})^2}{2\sigma_{i,k}^2}}$$

$$u_{i,k_{MLE}} = \frac{\sum_{x \in class_k} x_i}{|class_k|}$$

$$\sigma_{i,k_{MLE}}^{2} = \frac{\sum_{x \in class_{k}} \left(x_{i} - \mu_{i,k_{MLE}}\right)^{2}}{|class_{k}|}$$

 $class_k$: the set of examples of class k.

Recap: how to learn a Gaussian.

Unbiased:
$$\sigma_{i,k_{MLE}}^2 = \frac{\sum_{x \in class_k} (x_i - \mu_{i,k})^2}{|class_k| - 1}$$

Introduction to Machine Learning

Amo G. Tong

Summary

• Naïve Bayes Classifier

$$\operatorname{argmax} \Pr[v|x]$$

=argmax
$$Pr[x|v] Pr[v]$$

=argmax
$$\prod \Pr[x_i | v] \Pr[v]$$

Assume a form for them. Estimate it by data.

Introduction to Machine Learning

Amo G. Tong

Summary

- Baves optimal classifier
 - Review
- Naïve Bayes Classifier
 - · What is the assumption? Why we need it?
 - · When the attributes are discrete
 - How to estimate the parameters?
 - · m-estimate
 - · When the attributes are continuous
 - · Gaussian naïve Bayes.
 - · How to estimate the parameters?

	Sky	Water	Forecast	EnjoySport
)	Sunny	Warm	Same	Yes

Text classification

x= 'I have a red apple and a blue apple'

Introduction to Machine Learning

Amo G. Tong

Group Discussion

• A text classification example (from Dan Jurafsky)

	Doc	Words	Class
Training	1	Chinese Beijing Chinese	С
	2	Chinese Chinese Shanghai	С
	3	Chinese Macao	С
	4	Tokyo Japan Chinese	j
Test	5	Chinese Chinese Tokyo Japan	?

- We use m-estimate where (m=number of distinct word) and (p=1/m).
- m=6, p=1/6

Pr[j]=1/4

• $\Pr[j|x] = \left(\frac{1}{4}\right) * \left(\frac{2}{9}\right)^3 * \left(\frac{2}{9}\right) * \left(\frac{2}{9}\right) \approx 0.0001$ • $\Pr[c|x] = \left(\frac{3}{4}\right) * \left(\frac{3}{7}\right)^3 * \left(\frac{1}{14}\right) * \left(\frac{1}{14}\right) \approx 0.0003$

• y = (Japan, Macao, Macao)

Introduction to Machine Learning

Pr[c]=3/4

Pr[Chinese|j]=(1+1)/(3+6)=2/9 Pr[Tokyo|j]=(1+1)/(3+6)=2/9 Pr[Japan|j]=(1+1)/(3+6)=2/9

Pr[Chinese | c]=(5+1)/(8+6)=3/7 Pr[Tokyo|c]=(0+1)/(8+6)=1/14 Pr[Japan|c]=(0+1)/(8+6)=1/14

Amo G. Tong