

# CS 581 – ADVANCED ARTIFICIAL INTELLIGENCE

## TOPIC: LEARNING – NAÏVE BAYES



**Mustafa Bilgic**



<http://www.cs.iit.edu/~mbilgic>



<https://twitter.com/bilgicm>

# CLASSIFICATION

- AI → ML → Supervised Learning → Classification
- Given a dataset  $D = \{\langle X_i, y_i \rangle\}$  where
  - $X_i$  is the input
  - $y_i$  is the discrete-valued output
- Learn a function  $f(X_j) \rightarrow y_j$
- We would like  $f$  to **generalize** to **unseen** data
  - As opposed to memorize the given data

# CLASSIFICATION EXAMPLES

- Email classification
- Medical diagnosis
- Face recognition
- Optical character/digit recognition
- Sentiment classification
- ...

# ALGORITHMS

- Decision trees
- Nearest neighbor classification
- Naïve Bayes
- Logistic regression
- Support vector machines
- Neural networks
- ...

# GENERALIZATION

- The purpose of  $f$ 
  - Is not to memorize the “seen” data
  - Is to generalize to “unseen” data
- We need a performance metric
- We need to test  $f$ ’s performance on a dataset that it has not seen

# TYPES OF ERRORS – CLASSIFICATION

- Assume a target/positive class
  - Spam, HasHeartDisease, etc.
- *False positive*
  - Falsely classifying an object as positive
    - E.g., classifying a legitimate email as spam, diagnosing a healthy patient as having heart disease, and so on
  - Also called *Type I* error
- *False negative*
  - Falsely classifying an object as negative
    - E.g., classifying a spam email as not-spam, claiming that a heart-disease patient is healthy, and so on
  - Also called *Type II* error

# A FEW PERFORMANCE MEASURES

- 0/1 loss; error or accuracy
- Precision
- Recall
- F1
- Log-loss
- Domain-specific performance measures,
- ...

# CONFUSION MATRIX

		Predicted Class	
		Positive	Negative
Actual Class	Positive	True Positive	False Negative
	Negative	False Positive	True Negative



# ACCURACY

		Predicted Class	
		Positive	Negative
Actual Class	Positive	True Positive	False Negative
	Negative	False Positive	True Negative

$$Accuracy = \frac{Num\ Correct}{Data\ Size} = \frac{TP + TN}{TP + TN + FP + FN}$$

# PRECISION

		Predicted Class	
		Positive	Negative
Actual Class	Positive	True Positive	False Negative
	Negative	False Positive	True Negative

$$Precision = \frac{True\ Positive}{Predicted\ Positive} = \frac{TP}{TP + FP}$$

## TRUE POSITIVE RATE – RECALL – SENSITIVITY

		Predicted Class	
		Positive	Negative
Actual Class	Positive	True Positive	False Negative
	Negative	False Positive	True Negative

$$TPR = Recall = \frac{\text{True Positive}}{\text{Actual Positive}} = \frac{TP}{TP + FN}$$

## TRUE NEGATIVE RATE – SPECIFICITY

		Predicted Class	
		Positive	Negative
Actual Class	Positive	True Positive	False Negative
	Negative	False Positive	True Negative

$$TNR = Specificity = \frac{True\ Negative}{Actual\ Negative} = \frac{TN}{TN + FP}$$

## FALSE POSITIVE RATE – FALL-OUT

		Predicted Class	
		Positive	Negative
Actual Class	Positive	True Positive	False Negative
	Negative	False Positive	True Negative

$$FPR = FallOut = \frac{False\ Positive}{Actual\ Negative} = \frac{FP}{TN + FP}$$

## FALSE NEGATIVE RATE – MISS RATE

		Predicted Class	
		Positive	Negative
Actual Class	Positive	True Positive	False Negative
	Negative	False Positive	True Negative

$$FNR = Miss Rate = \frac{False\ Negative}{Actual\ Positive} = \frac{FN}{TP + FN}$$

# F1

		Predicted Class	
		Positive	Negative
Actual Class	Positive	True Positive	False Negative
	Negative	False Positive	True Negative

$$F1 = \frac{2 * Precision * Recall}{Precision + Recall}$$

# SPLITTING THE DATASET

1. Train-test splits
2. Train-validation-test splits
3. Cross-validation



# TRAIN-TEST SPLIT

- Randomly split the data into two disjoint sets
- A typical approach:  $2/3$  for train and  $1/3$  for test
- Train your model on training data and evaluate it on the test data
  - Use your favorite performance metric
- Report your performance as the expected performance on unseen data
- Caveats:
  - You need a large dataset for this to work
  - You cannot tune your parameters on the test data

# TRAIN-VALIDATION-TEST SPLIT

- Split your data into three disjoint sets
  - Train, validation, test
- Train your model(s) on the training data
- Evaluate your model(s) on the validation data
- Pick the model that performs best on the validation data
- Test the model on the test data, and report its performance
- Caveat:
  - You need a really big dataset for this to work

# CROSS-VALIDATION

- Split your data into  $k$  disjoint sets
- Each time, one set is the test set and the rest is the training set
- See OneNote for more detailed explanation and illustration

# REAL LIFE MEASURES

- Not as clean as the ones we discussed
- Imagine self-driving cars, medical diagnosis, crime prediction, fraud detection, and so on
- Usually, there is not a single performance measure
- Performance is handled on a case-by-case basis; not on an aggregate level

# NAÏVE BAYES

# BAYES CLASSIFIER

- Input:  $\vec{X} = \langle X_1, X_2, \dots, X_n \rangle$
- Output:  $Y$
- Bayes classifier

$$P(Y | \vec{X}) = \frac{P(\vec{X} | Y)P(Y)}{P(\vec{X})} = \frac{P(Y)P(X_1, X_2, \dots, X_n | Y)}{P(X_1, X_2, \dots, X_n)}$$

$$P(X_1, X_2, \dots, X_n) = \sum_y P(Y = y)P(X_1, X_2, \dots, X_n | Y = y)$$

Assuming all variables are binary, how many independent parameters are needed for the Bayes classifier?

# BAYES CLASSIFIER

- Assume a binary classification task, where the label  $Y$  is *spam* or  $\sim\text{spam}$
- Assume  $P(\text{spam}) = 0.4$
- If you have seen an email  $X$   $a$  times as *spam* and  $b$  times as  $\sim\text{spam}$ 
  - Using an MLE estimate for  $P(X|Y)$ , what is  $P(Y|X)$ ?
  - Using an LS estimate for  $P(X|Y)$ , what is  $P(Y|X)$ ?
  - What happens when either or both of  $a$  and  $b$  are zero?

# NAÏVE BAYES ASSUMPTION

$$X_i \perp X_j \mid Y$$



# NAÏVE BAYES

Bayes rule:

$$P(Y | X_1, X_2, \dots, X_n) = \frac{P(Y)P(X_1, X_2, \dots, X_n | Y)}{\sum_y P(y)P(X_1, X_2, \dots, X_n | y)}$$

Assuming  $X_i \perp X_j | Y$ ,  
naïve Bayes:

$$P(Y | X_1, X_2, \dots, X_n) = \frac{P(Y) \prod P(X_i | Y)}{\sum_y P(y) \prod P(X_i | y)}$$

Assuming all variables are binary, how many independent parameters are needed for the naive Bayes classifier?

# EXAMPLE

- See OneNote

# NAÏVE BAYES IMPLEMENTATIONS

- Bernoulli / categorical naïve Bayes
  - Features are assumed to be binary / categorical
- Multinomial naïve Bayes
  - $P(\vec{X} \mid y)$  is a multinomial distribution
- Gaussian naïve Bayes
  - Each  $p(x_i \mid y)$  is a Gaussian distribution

# READING

- <http://www.cs.cmu.edu/~tom/mlbook/NBayesLogReg.pdf>
- [https://en.wikipedia.org/wiki/Naive\\_Bayes\\_classifier](https://en.wikipedia.org/wiki/Naive_Bayes_classifier)
- [https://scikit-learn.org/stable/modules/naive\\_bayes.html](https://scikit-learn.org/stable/modules/naive_bayes.html)

# ZERO PROBABILITIES

- Assume feature  $X_i$  is T for a particular object.  
Further,
  - Assume  $P(X_i = T|yes) = 0$  and  $P(X_i = T|no) > 0$  and  $P(X_j | yes) > 0$  and  $P(X_j | no) > 0$  for all other features
    - What is  $P(yes | \vec{X})$ ?
  - Assume  $P(X_i = T|yes) = 0$  and  $P(X_i = T|no) = 0$  and  $P(X_j | yes) > 0$  and  $P(X_j | no) > 0$  for all other features
    - What is  $P(yes | \vec{X})$ ?
- One solution: use LS for the parameter estimates

# MULTIPLYING SEVERAL PROBABILITY NUMBERS

- Assume we have 10,000 features
- What is  $0.9^{10,000}$  using a computer?
- Try `math.pow(0.9, 10000)` in Python
- In Naïve Bayes,
  - $a = P(Y = T) \prod P(X_i|Y = T)$
  - $b = P(Y = F) \prod P(X_i|Y = F)$
  - $P(Y = T|\vec{X}) = \frac{a}{a+b}$
  - If  $a = b = 0$  in your code, then what?