

SLIDE-1

Data Mining Supervised:

- KNN
- Naïve Bayes
- Decision trees
- Random forest
- Neural networks
- Support vector machines

SLIDE-2

k-Nearest Neighbor Classifiers

SLIDE-3

1-Nearest Neighbor Classifier

Training Examples (Instances)
Some for each CLASS

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

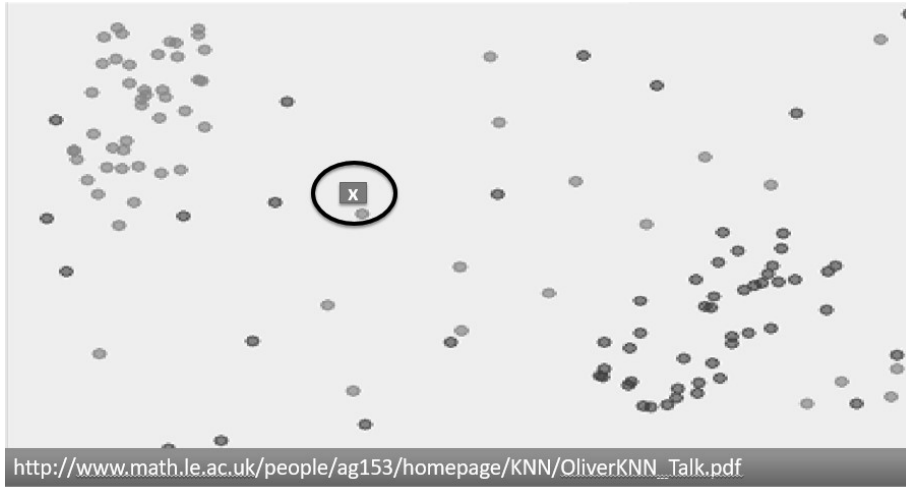
Test Examples
(What class to assign this?)

4

The diagram shows a 1-Nearest Neighbor Classifier. Training examples are organized into a grid with two columns. The left column contains four groups of four handwritten digits each, labeled 1, 3, 5, and 7. The right column contains three groups of three handwritten digits each, labeled 2, 4, 6, 8, and 0. The digit '4' in the second row of the right column is circled. To the right of the grid, a test example '4' is shown with an arrow pointing to the circled '4' in the training set. The text 'Test Examples (What class to assign this?)' is positioned above the test example.

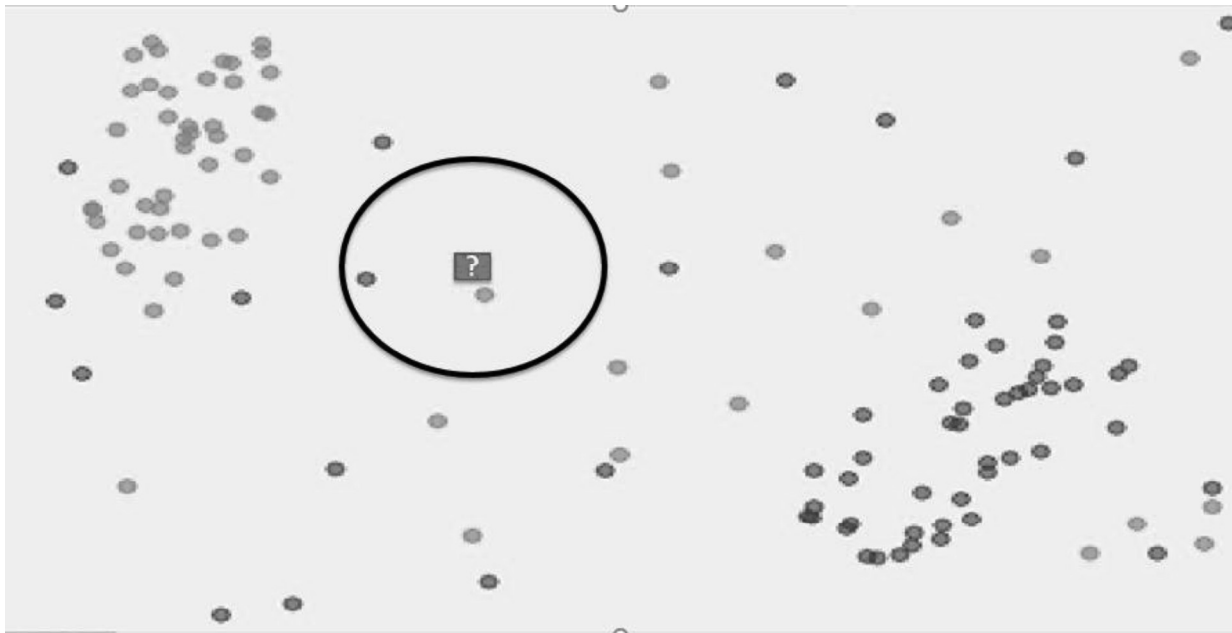
SLIDE-4

1-Nearest Neighbor



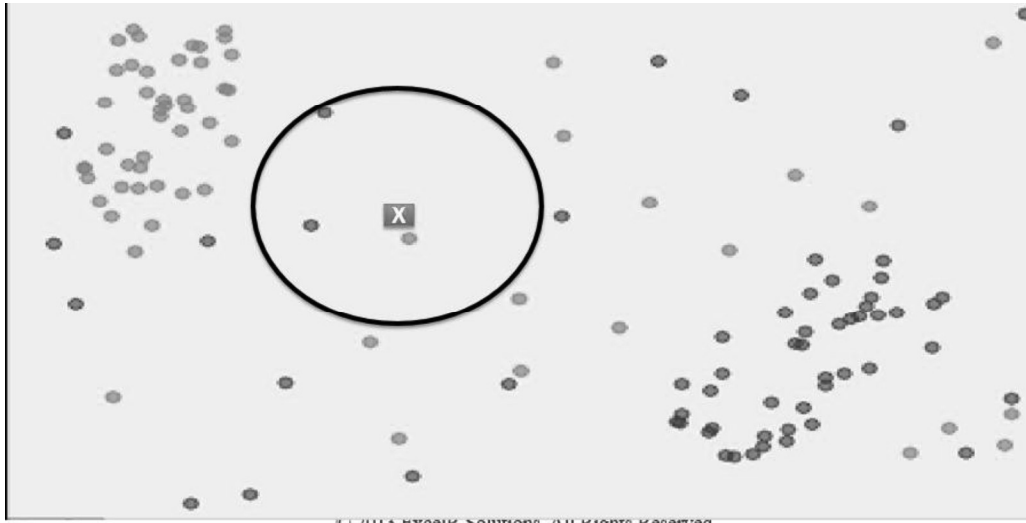
SLIDE-5

2-Nearest Neighbor



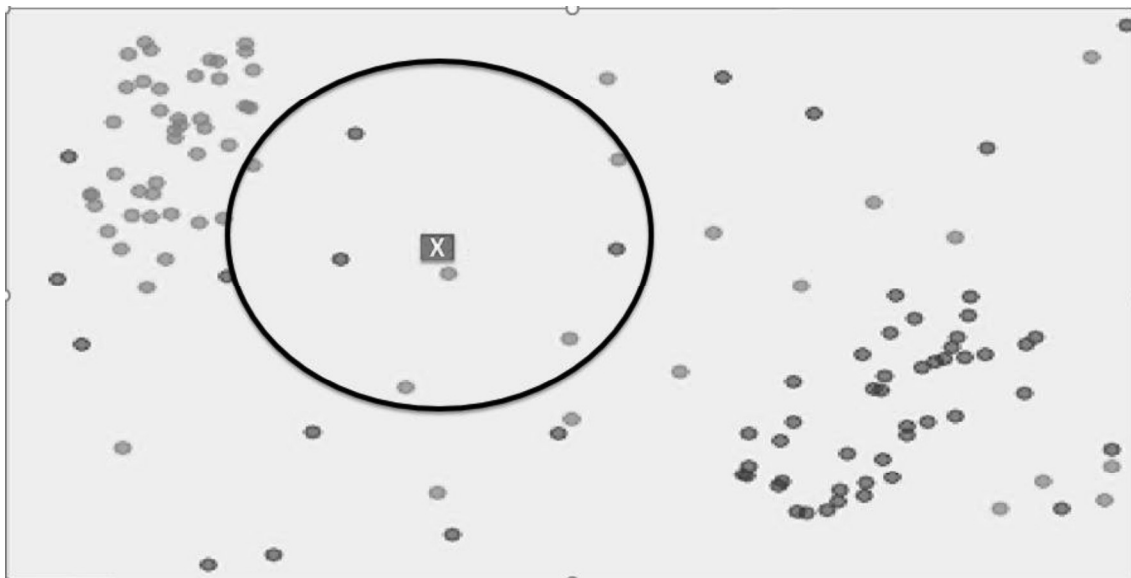
SLIDE-6

3-Nearest Neighbor



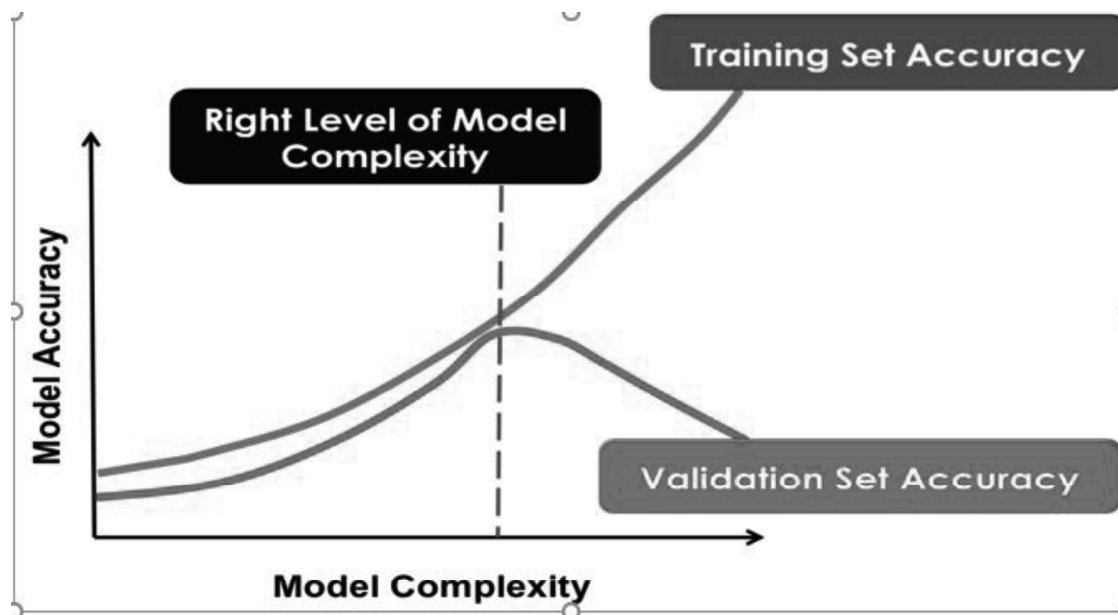
SLIDE-7

8-Nearest Neighbor



SLIDE-8

Controlling COMPLEXITY in k-NN



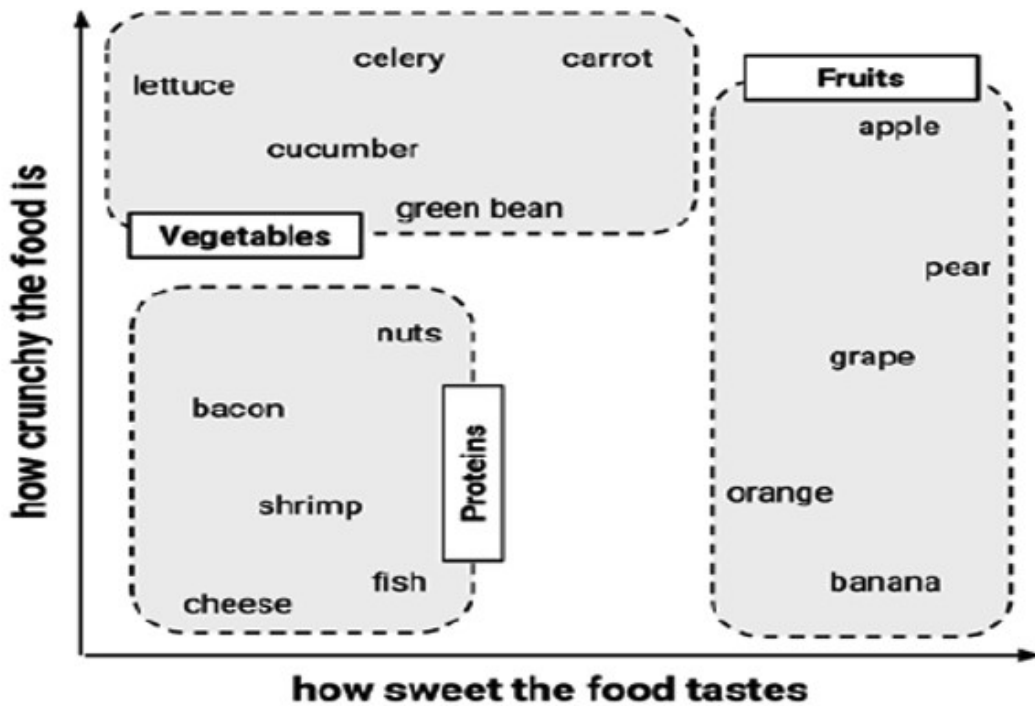
SLIDE-9

Ingredient	Sweetness	Crunchiness	Food type
apple	10	9	fruit
Bacon	1	4	protein
banana	10	1	fruit
carrot	7	10	vegetable
celery	3	10	vegetable
cheese	1	1	protein

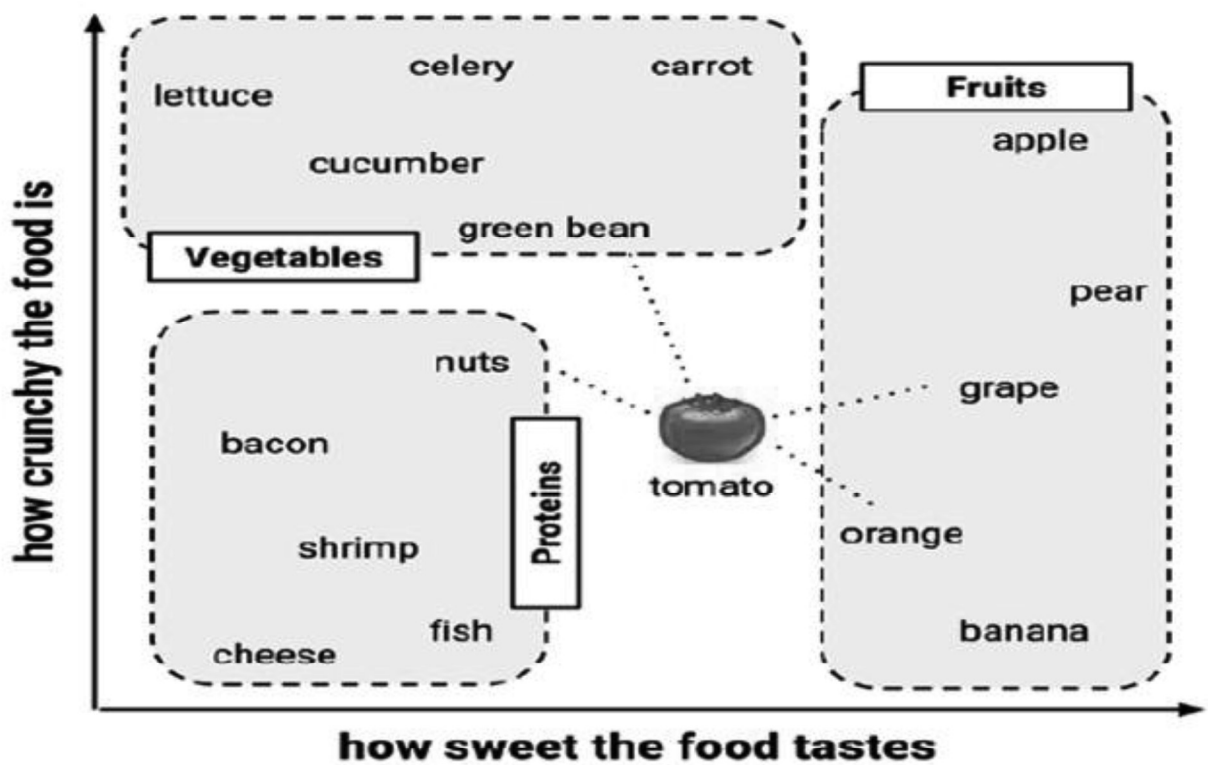
SLIDE-10



SLIDE-11



SLIDE-12



SLIDE-13

Measuring similarity with distance

$$\text{dist}(p, q) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_n - q_n)^2}$$

Which Class Tomoto belongs to given the feature values:

Tomato (*sweetness* = 6, *crunchiness* = 4),

SLIDE-14

Bayesian Classifiers

SLIDE-15

Understanding probability

The probability of an event is estimated from the observed data by dividing the number of trials in which the event occurred by the total number of trials

For instance, if it rained 3 out of 10 days with similar conditions as today, the probability of rain today can be estimated as $3 / 10 = 0.30$ or 30 percent.

Similarly, if 10 out of 50 prior email messages were spam, then the probability of any incoming message being spam can be estimated as $10 / 50 = 0.20$ or 20 percent.

For example, given the value $P(\text{spam}) = 0.20$, we can calculate $P(\text{ham}) = 1 - 0.20 = 0.80$

Note: The probability of all the possible outcomes of a trial must always sum to 1

SLIDE-16

Understanding probability cont..

For example, given the value $P(\text{spam}) = 0.20$, we can calculate $P(\text{ham}) = 1 - 0.20 = 0.80$

Because an event cannot simultaneously happen and not happen, an event is always mutually exclusive and exhaustive with its complement

The complement of event A is typically denoted A^c or A' .

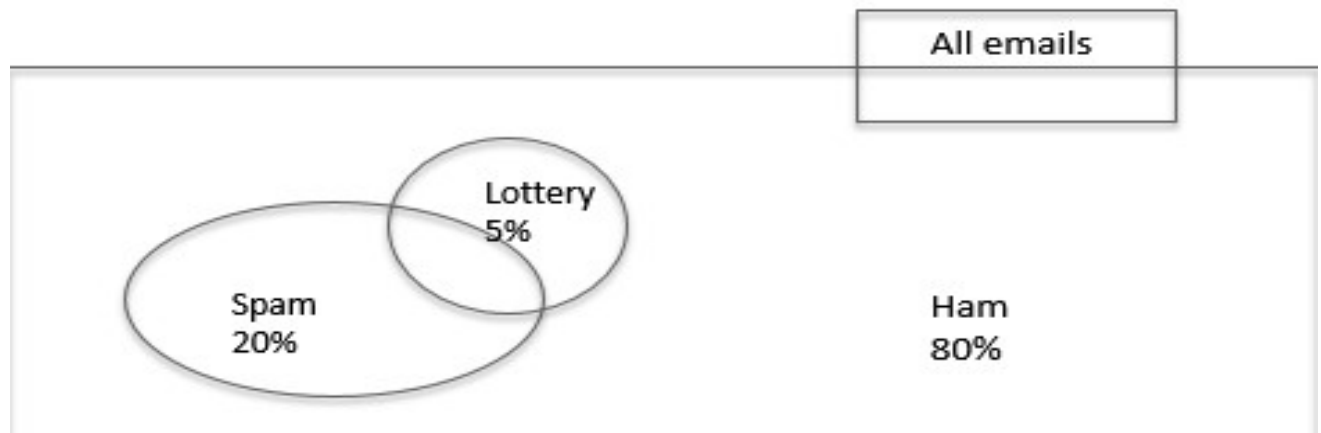
Additionally, the shorthand notation $P(\neg A)$ can be used to denote the probability of event A not occurring, as in $P(\neg \text{spam}) = 0.80$. This notation is equivalent to $P(A^c)$.



SLIDE-17

Understanding joint probability

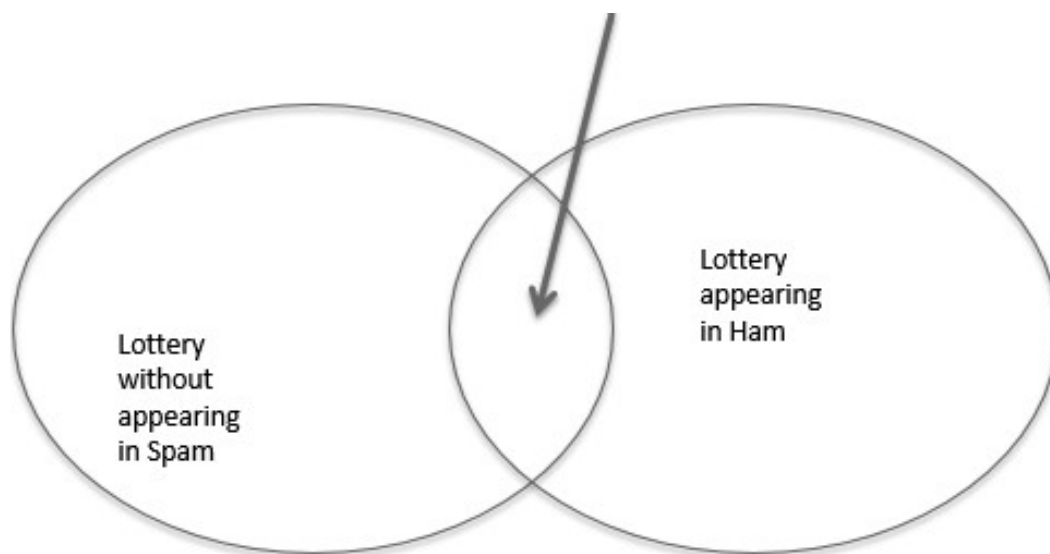
Often, we are interested in monitoring several nonmutually exclusive events for the same trial



SLIDE-18

Understanding joint probability

Lottery appearing in spam



Estimate the probability that both $P(\text{spam})$ and $P(\text{Spam})$ occur, which can be written as $P(\text{spam} \cap \text{Lottery})$. the notation $A \cap B$ refers to the event in which both A and B occur.

SLIDE-19

Calculating $P(\text{spam} \cap \text{Lottery})$ depends on the joint probability of the two events or how the probability of one event is related to the probability of the other.

If the two events are totally unrelated, they are called independent events

If $P(\text{spam})$ and $P(\text{Lottery})$ were independent, we could easily calculate $P(\text{spam} \cap \text{Lottery})$, the probability of both events happening at the same time.

Because 20 percent of all the messages are spam, and 5 percent of all the e-mails contain the word Lottery, we could assume that 1 percent of all messages are spam with the term Lottery.

More generally, for independent events A and B, the probability of both happening can be expressed as $P(A \cap B) = P(A) * P(B)$.

$$0.05 * 0.20 = 0.01$$

SLIDE-20

Bayes Rule

$$P(A|B) = P(A \cap B)/P(B) = P(B/A)P(A)/P(B)$$

- Bayes Rule: The most important Equation in ML!!

$$P(\text{CLASS}/\text{DATA}) = P(\text{CLASS}) P(\text{DATA}/\text{CLASS})/P(\text{DATA})$$

LHS → POSTERIOR PROBABILITY (Probability of class AFTER seeing the data)

RHS → DATA PRIOR (MARGINAL)

SLIDE-21

Naïve Bayes Classifier

SLIDE-22

Conditional Independence

$$P(\text{Fever}, \text{BodyAche} | \text{Viral}) = P(\text{Fever} | \text{Viral}) * P(\text{BodyAche} | \text{Viral})$$

Viral Infection

- Fever
- Body Ache

$$P(\text{Fever}, \text{BodyAche}) \neq P(\text{Fever})P(\text{BodyAche})$$

- Simple Independence between two variables:

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

- Class Conditional Independence assumption:

$$P(X_1, X_2 / C) = P(X_1 / C)P(X_2 / C)$$

SLIDE-23

Naïve Bayes Classifier

Conditional Independence among variables given Classes!

$$P(C|X_1, X_2, \dots, X_D) = \frac{P(C)P(X_1, X_2, \dots, X_D|C)}{\sum_{C'} P(X_1, X_2, \dots, X_D|C')} = \frac{P(C) \prod_{d=1}^D P(X_d|C)}{\sum_{C'} \prod_{d=1}^D P(X_d|C')}$$

- Simplifying assumption
- Baseline model especially when large number of features
- Taking log and ignoring denominator:

$$\log(P(C|X_1, X_2, \dots, X_D)) \propto \log(P(C)) + \sum_{d=1}^D \log(P(X_d|C))$$

SLIDE-24

Naïve Bayes Classifier for Categorical Valued Variables

Let's Naïve Bayes!

$$\log(P(C|X_1, X_2, \dots, X_D)) \propto \log(P(C)) + \sum_{d=1}^D \log(P(X_d|C))$$

Class Prior Parameters:

$$P(\text{Like} = Y) = ???$$

$$P(\text{Like} = N) = ???$$

Class Conditional Likelihoods

$$P(\text{Color} = \text{Red} | \text{Like} = Y) = ????$$

$$P(\text{Color} = \text{Red} | \text{Like} = N) = ????$$

...

$$P(\text{Shape} = \text{Triangle} | \text{Like} = N) = ????$$

#EXMPLS	COLOR	SHAPE	LIKE
20	Red	Square	Y
10	Red	Circle	Y
10	Red	Triangle	N
10	Green	Square	N
5	Green	Circle	Y
5	Green	Triangle	N
10	Blue	Square	N
10	Blue	Circle	N
20	Blue	Triangle	Y

SLIDE-26**Parameter Estimation**

$$\log(P(C | X_1, X_2, \dots, X_D)) \propto \log(P(C)) + \sum_{d=1}^D \log(P(X_d | C))$$

- What / How many Parameters?

P(C)

$$= \frac{N(c)}{N} \approx \frac{N(c) + \lambda}{N(c) + \lambda \times |\mathbf{C}|}$$

- **Class Priors:**
- **Conditional Probabilities:**

$$X_d = \{v_1^{(d)}, v_2^{(d)}, \dots, v_{M_d}^{(d)}\}$$

$$P(v_m^{(d)} | c) = \frac{N(v_m^{(d)}, c)}{N(c)} \cup \frac{N(v_m^{(d)}, c) + \lambda}{N(c) + M_d \lambda}$$

SLIDE-27

Naïve Bayes Classifier for Text Classifier

Text Classification Example

- Doc1 = {buy two shirts get one shirt half off}
- Doc2 = {get a free watch. send your contact details now}
- Doc3 = {your flight to chennai is delayed by two hours}
- Doc4 = {you have three tweets from @sachin}

Four Class Problem:

- Spam, $P(\text{promo}|\text{doc1}) = 0.84$
- Promotions, $P(\text{spam}|\text{doc2}) = 0.94$
- Social, $P(\text{main}|\text{doc3}) = 0.75$
- Main $P(\text{social}|\text{doc4}) = 0.91$

SLIDE-29

Bag-of-Words Representation

- Structured (e.g. Multivariate) data – fixed number of features
- Unstructured (e.g. Text) data
 - Arbitrary length documents,
 - High dimensional feature space (many words in vocabulary),
 - Sparse (small fraction of vocabulary words present in a doc.)
- Bag-of-Words Representation:
 - Ignore Sequential order of words
 - Represent as a Weighted-Set – Term Frequency of each term
- RawDoc = {buy two shirts get one shirt half off}

- Stemming = {buy two shirt get one shirt half off}
- BoW's = {buy:1, two:1, shirt:2, get:1, one:1, half:1, off:1}

SLIDE-30

Naïve Bayes Classifier with BoW

BoW = {buy:1, two:1, shirt:2, get:1, one:1, half:1, off:1}

- Make an “independence assumption” about words | class

$P(\text{doc1}|\text{promo}) =$

$P(\text{buy:1,two:1,shirt:2,get:1,one:1,half:1,off:1}|\text{promo})$

$= P(\text{buy:1}|\text{promo}) * P(\text{two:1}|\text{promo}) * P(\text{shirt:2}|\text{promo}) * P(\text{get:1}|\text{promo}) * P(\text{one:1}|\text{promo}) * P(\text{free:1}|\text{promo})$

$= P(\text{buy}|\text{promo})^1 * P(\text{two}|\text{promo})^1 * P(\text{shirt}|\text{promo})^2$
 $* P(\text{get}|\text{promo})^1 * P(\text{one}|\text{promo})^1 * P(\text{free}|\text{promo})^1$

SLIDE-31

Naïve Bayes Text Classifiers

- Log Likelihood of document given class.

$$doc = \{tf(w_m)\}_{m=1}^M$$

$tf(w_m)$ = Number of times word w_m occurs in doc

$$P(doc|class) = P(w_1|class)^{tf(w_1)} P(w_2|class)^{tf(w_2)} \dots P(w_M|class)^{tf(w_M)}$$

- Parameters in Naïve Bayes Text classifiers:

$P(w_m|c)$ = Probability that word w_m occurs in documents of class c

$P(shirt|promo), P(free|spam), P(buy|spam), P(buy|promo), \dots$

Number of parameters = ??

SLIDE-32

Naïve Bayes Parameters

- Likelihood of a word given class. For each word, each class.
- Estimating these parameters from data:

$P(w_m|c)$ = Probability that word w_m occurs in documents of class c

- Estimating these parameters from data:

$N(w_m, c)$ = Number of times word w_m occurs in documents of class c

$$N(\text{free}, \text{spam}) = \sum_{\text{doc} \in \text{spam}} \text{tf}(\text{free} | \text{doc})$$

$$N(\text{free}, \text{promo}) = \sum_{\text{doc} \in \text{promo}} \text{tf}(\text{free} | \text{doc})$$

SLIDE-33

- Likelihood of a word given class. For each word, each class.

$P(w_m | c)$ = Probability that word w_m occurs in documents of class c

- Normalize these counts to probabilities.

$$P(\text{free} | \text{spam}) = \frac{N(\text{free}, \text{spam})}{\sum_{m'=1}^M N(w_{m'}, \text{spam})} \approx \frac{1 + N(\text{free}, \text{spam})}{M + \sum_{m'=1}^M N(w_{m'}, \text{spam})}$$

- Smoothing is done to make sure we don't get zero products

SLIDE-34

Bayesian Classifier Multi-variate real-valued data

SLIDE-35

Bayes Rule

Class Prior

Data Likelihood given Class

$$P(\text{Class}|\text{Data}) = \frac{P(\text{Class})P(\text{Data}|\text{Class})}{P(\text{Data})}$$

Posterior Probability

Data Prior (Marginal)

(Probability of class AFTER seeing the data)

$$P(c|\mathbf{x}) = \frac{P(c)P(\mathbf{x}|c)}{P(\mathbf{x})}$$

$$\mathbf{x} \in R^D$$

SLIDE-36

Simple Bayesian Classifier

$$P(c|\mathbf{x}) = \frac{P(c)P(\mathbf{x}|c)}{P(\mathbf{x})}$$

$P(\mathbf{x}|c) \rightarrow$ Each Class Conditional Probability is assumed to be a Uni-Modal (Single Cloud) (NORMAL) Distribution.

$$P(\mathbf{x}|c) = N(\mathbf{x}|\mathbf{m}_c, \Sigma_c) = \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma_c|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{m}_c)^T \Sigma_c^{-1} (\mathbf{x} - \mathbf{m}_c)\right)$$

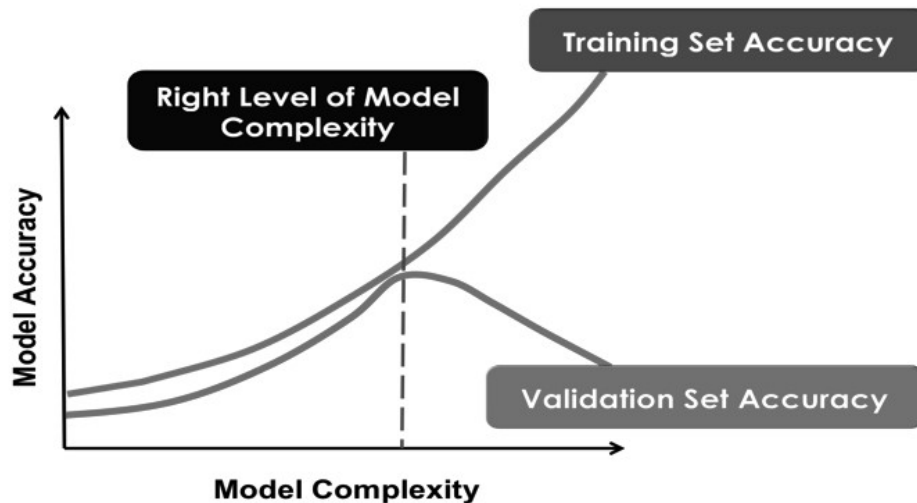
$$\text{Sum: } \int_{\mathbf{x} \in R^D} P(\mathbf{x}|c) d\mathbf{x} = 1$$

$$\text{Mean: } \int_{\mathbf{x} \in R^D} \mathbf{x} P(\mathbf{x}|c) d\mathbf{x} = \mathbf{m}_c$$

$$\text{Co-Variance: } \int_{\mathbf{x} \in R^D} (\mathbf{x} - \mathbf{m}_c)(\mathbf{x} - \mathbf{m}_c)^T P(\mathbf{x}|c) d\mathbf{x} = \Sigma_c$$

SLIDE-37

Controlling COMPLEXITY



SLIDE-38

Spherical-SAME Covariance/Class

$$P(\mathbf{x}|c) = \frac{1}{Z(\mathbf{m}_c, \Sigma_c)} \exp\left(-\frac{1}{2} \Delta(\mathbf{x}, \mathbf{m}_c | \Sigma_c)^2\right)$$

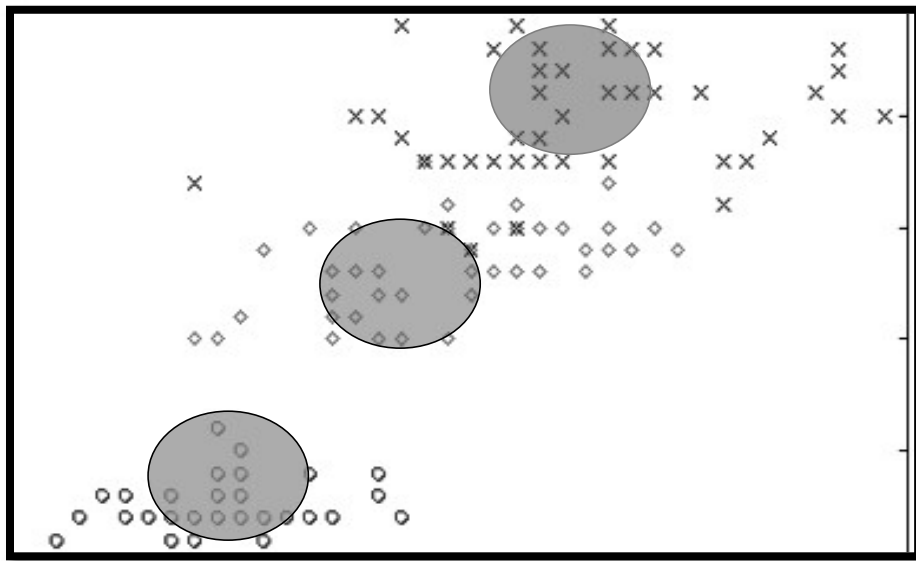
$$= \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma_c|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\mathbf{x} - \mathbf{m}_c)^T \Sigma_c^{-1} (\mathbf{x} - \mathbf{m}_c)\right)$$

$$\Sigma_c = \begin{bmatrix} \sigma & 0 & 0 & 0 \\ 0 & \sigma & 0 & 0 \\ 0 & 0 & \sigma & 0 \\ 0 & 0 & 0 & \sigma \end{bmatrix}$$

$$= \frac{1}{(\sqrt{2\pi})^D \sigma} \exp\left(-\frac{1}{2} \sum_{d=1}^D \left(\frac{x_d - m_{c,d}}{\sigma}\right)^2\right)$$

Euclidian Distance

$$\Delta(\mathbf{x}, \mathbf{m}_c | \Sigma_c) = \sqrt{\sum_{d=1}^D \left(\frac{x_d - m_{c,d}}{\sigma}\right)^2}$$



SLIDE-39

Spherical-**DIFFERENT** Covariance/Class

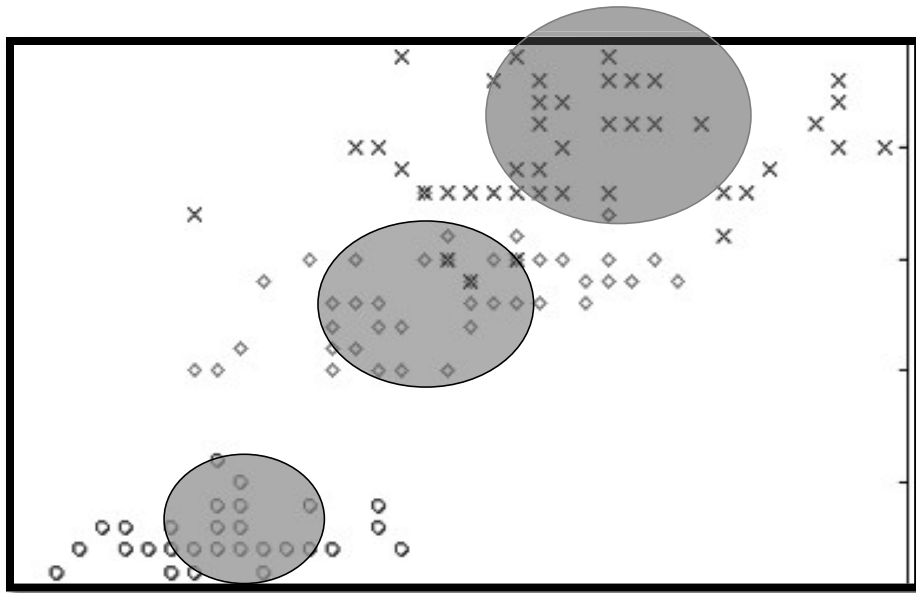
$$\begin{aligned}
 P(\mathbf{x}|c) &= \frac{1}{Z(\mathbf{m}_c, \Sigma_c)} \exp\left(-\frac{1}{2} \Delta(\mathbf{x}, \mathbf{m}_c | \Sigma_c)^2\right) \\
 &= \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma_c|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\mathbf{x} - \mathbf{m}_c)^T \Sigma_c^{-1} (\mathbf{x} - \mathbf{m}_c)\right)
 \end{aligned}$$

$$\Sigma_c = \begin{bmatrix} \sigma_c & 0 & 0 & 0 \\ 0 & \sigma_c & 0 & 0 \\ 0 & 0 & \sigma_c & 0 \\ 0 & 0 & 0 & \sigma_c \end{bmatrix}$$

$$= \frac{1}{(\sqrt{2\pi})^D \sigma} \exp\left(-\frac{1}{2} \sum_{d=1}^D \left(\frac{x_d - m_{c,d}}{\sigma}\right)^2\right)$$

Class Scaled Euclidian Distance

$$\Delta(\mathbf{x}, \mathbf{m}_c | \Sigma_c) = \sqrt{\sum_{d=1}^D \left(\frac{x_d - m_{c,d}}{\sigma_c}\right)^2}$$



SLIDE-40

Independent-Covariance/Class

$$P(\mathbf{x}|c) = \frac{1}{Z(\mathbf{m}_c, \Sigma_c)} \exp\left(-\frac{1}{2} \Delta(\mathbf{x}, \mathbf{m}_c | \Sigma_c)^2\right)$$

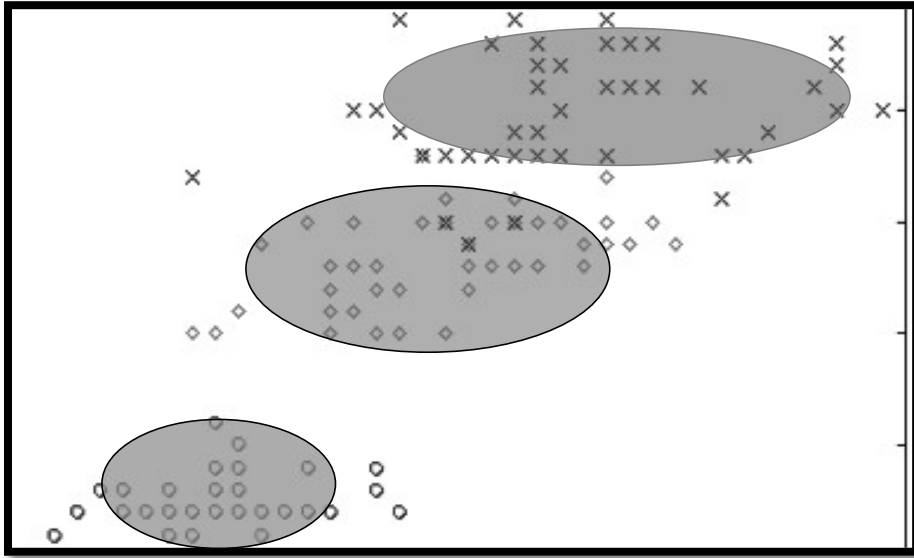
$$= \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma_c|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\mathbf{x} - \mathbf{m}_c)^T \Sigma_c^{-1} (\mathbf{x} - \mathbf{m}_c)\right)$$

$$\Sigma_c = \begin{bmatrix} \sigma_{c,1} & 0 & 0 & 0 \\ 0 & \sigma_{c,2} & 0 & 0 \\ 0 & 0 & \sigma_{c,3} & 0 \\ 0 & 0 & 0 & \sigma_{c,4} \end{bmatrix}$$

$$= \frac{1}{(\sqrt{2\pi})^D \sigma} \exp\left(-\frac{1}{2} \sum_{d=1}^D \left(\frac{x_d - m_{c,d}}{\sigma}\right)^2\right)$$

Class/Dim. Scaled Euclidean Distance

$$\Delta(\mathbf{x}, \mathbf{m}_c | \Sigma_c) = \sqrt{\sum_{d=1}^D \left(\frac{x_d - m_{c,d}}{\sigma_{c,d}} \right)^2}$$



SLIDE-41

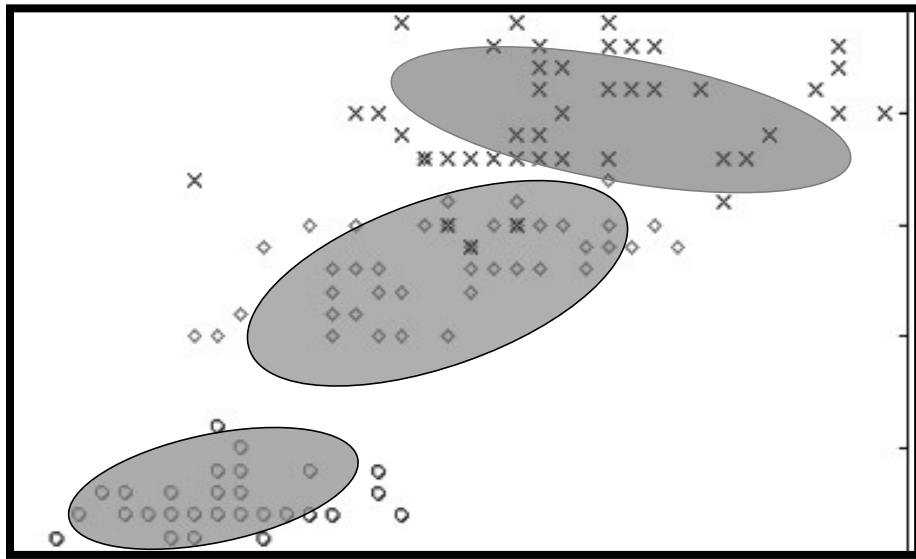
FULL-Covariance/Class

$$\begin{aligned} P(\mathbf{x}|c) &= \frac{1}{Z(\mathbf{m}_c, \Sigma_c)} \exp\left(-\frac{1}{2} \Delta(\mathbf{x}, \mathbf{m}_c | \Sigma_c)^2\right) \\ &= \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma_c|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\mathbf{x} - \mathbf{m}_c)^T \Sigma_c^{-1} (\mathbf{x} - \mathbf{m}_c)\right) \end{aligned}$$

$$\Sigma_c = \begin{bmatrix} \sigma_{1,1}^{(c)} & \sigma_{1,2}^{(c)} & \sigma_{1,3}^{(c)} & \sigma_{1,4}^{(c)} \\ \sigma_{2,1}^{(c)} & \sigma_{2,2}^{(c)} & \sigma_{2,3}^{(c)} & \sigma_{2,4}^{(c)} \\ \sigma_{3,1}^{(c)} & \sigma_{3,2}^{(c)} & \sigma_{3,3}^{(c)} & \sigma_{3,4}^{(c)} \\ \sigma_{4,1}^{(c)} & \sigma_{4,2}^{(c)} & \sigma_{4,3}^{(c)} & \sigma_{4,4}^{(c)} \end{bmatrix}$$

Mahalanobis Distance

$$\Delta(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T \Sigma^{-1} (\mathbf{x} - \mathbf{y})}$$



SLIDE-42

Summary

- Classifier = Partitions Input Space into “Pure” regions
- Rule Based Classifiers:
 - Version Space

- Decision Trees
- Descriptive Classification:
 - Bayesian Classifiers
- Discriminative Classification:
 - Perceptron
 - Logistic Regression