CS217: Artificial Intelligence and Machine Learning (associated lab: CS240)

Pushpak Bhattacharyya, Nihar Ranjan Sahoo CSE Dept., IIT Bombay

Week8 of 3mar25, Prolog, Midsem discussion, SVM primal, dual, Kernel Trick

Main points covered: week7 of 17feb25

Inferencing in Predicate Calculus

- Forward chaining
 - Given P, $P \rightarrow Q$, to infer Q
 - P, match L.H.S of
 - Assert Q from *R*.*H*.*S*
- Backward chaining
 - Q, Match R.H.S of $P \rightarrow Q$
 - assert P
 - Check if P exists
- Resolution Refutation
 - Negate goal
 - Convert all pieces of knowledge into clausal form (disjunction of literals)
 - See if contradiction indicated by null clause ____ can be derived

Wh-Questions and Knowledge



Knowledge Representation of Complex Sentence

"In every city there is a thief who is beaten by every policeman in the city"

 $\forall x [city(x) \rightarrow \{ \exists y ((thief(y) \land lives_{in}(y, x)) \land \forall z (policeman(z, x) \rightarrow beaten_{by}(z, y))) \}]$

Interpretation in Logic

- Logical expressions or formulae are "FORMS" (placeholders) for whom <u>contents</u> are created through interpretation.
- Example:

$$\exists F[\{F(a) = b\} \land \forall x \{P(x) \to (F(x) = g(x, F(h(x))))\}]$$

- This is a Second Order Predicate Calculus formula.
- Quantification on 'F' which is a function.

Examples

- Interpretation:1 D=N (natural numbers) a = 0 and b = 1 $x \in N$ P(x) stands for x > 0q(m,n) stands for $(m \times n)$ h(x) stands for (x - 1)
- Above interpretation defines Factorial

Examples (contd.)

• Interpretation:2 $D=\{\text{strings}\}$ $a = b = \lambda$ P(x) stands for "x is

P(x) stands for "x is a non empty string" g(m, n) stands for "append head of m to n" h(x) stands for tail(x)

 Above interpretation defines "reversing a string"

End main points



Introduction

- PROgramming in LOGic
- Emphasis on what rather than how

Problem in Declarative Form

Logic Machine

Basic Machine

A Typical Prolog program

Compute_length ([],0). Compute_length ([Head|Tail], Length):-Compute_length (Tail,Tail_length), Length is Tail_length+1. High level explanation:

The length of a list is 1 plus the length of the tail of the list, obtained by removing the first element of the list.

This is a declarative description of the computation.

Fundamentals

(absolute basics for writing Prolog Programs)

Facts

- John likes Mary
 - like(john,mary)
- Names of relationship and objects must begin with a lower-case letter.
- Relationship is written *first* (typically the *predicate* of the sentence).
- Objects are written separated by commas and are enclosed by a pair of round brackets.
- The full stop character `.' must come at the end of a fact.



| Predicate | Interpretation |
|------------------------|--------------------------------|
| valuable(gold) | Gold is valuable. |
| owns(john,gold) | John owns gold. |
| father(john,mary) | John is the father of Mary |
| gives (john,book,mary) | John gives the book to Mary |

Questions

- Questions based on facts
- Answered by matching

Two facts *match* if their predicates are same (spelt the same way) and the arguments each are same.

- If matched, prolog answers *yes*, else *no*.
- *No* does not mean falsity.

Prolog does theorem proving

- When a question is asked, prolog tries to match *transitively*.
- When no match is found, answer is *no*.
- This means *not provable* from the given facts.

Variables

- Always begin with a capital letter
 - ?- likes (john,X).
 - ?- likes (john, Something).
- But not
 - ?- likes (john, something)

Example of usage of variable

Facts: likes(john,flowers). likes(john,mary). likes(paul,mary). Question: ?- likes(john,X) Answer: X=flowers and wait ; mary ; no

Conjunctions

- Use `,' and pronounce it as and.
- Example
 - Facts:
 - likes(mary,food).
 - likes(mary,tea).
 - likes(john,tea).
 - likes(john,mary)
- _ ?-
- likes(mary,X),likes(john,X).
- Meaning is anything liked by Mary also liked by John?

Backtracking (an inherent property of prolog programming)

likes(mary,X),likes(john,X)

-likes(mary,food) likes(mary,tea) likes(john,tea) likes(john,mary)

First goal succeeds. X=food
 Satisfy *likes(john,food)*

Backtracking (continued)

Returning to a marked place and trying to resatisfy is called *Backtracking*



Backtracking (continued)



First goal succeeds again, *X=tea* Attempt to satisfy the *likes(john,tea)*

Backtracking (continued)



Second goal also succeeds
 Prolog notifies success and waits for a reply

Rules

- Statements about *objects* and their relationships
- Expess
 - If-then conditions
 - I use an umbrella if there is a rain
 - use(i, umbrella) :- occur(rain).
 - Generalizations
 - All men are mortal
 - *mortal(X) :- man(X).*
 - Definitions
 - An animal is a bird if it has feathers
 - bird(X) :- animal(X), has_feather(X).

Syntax

- Read `:-' as `if'.
- E.G.
 - likes(john,X) :- likes(X,cricket).
 - "John likes X if X likes cricket".
 - *i.e.,* "John likes anyone who likes cricket".
- Rules always end with `.'.

Another Example

sister_of (X,Y):- female (X), parents (X, M, F), parents (Y, M, F).

X is a sister of Y is X is a female and X and Y have same parents

Question Answering in presence of *rules*

- Facts
 - male (ram).
 - male (shyam).
 - female (sita).
 - female (gita).
 - parents (shyam, gita, ram).
 - parents (sita, gita, ram).

Question Answering: Y/N type: *is sita the sister of shyam?*



Question Answering: wh-type: whose sister is sita?



Rules

- Statements about *objects* and their relationships
- Express
 - If-then conditions
 - I use an umbrella if there is a rain
 - use(i, umbrella) :- occur(rain).
 - Generalizations
 - All men are mortal
 - mortal(X) :- man(X).
 - Definitions
 - An animal is a bird if it has feathers
 - bird(X) :- animal(X), has_feather(X).

Support Vector Machine (SVM)

Introduction

- Support Vector Machine (SVM): Learns a linear separator for separating instances belonging to two different classes
 - In case of 1 dimensional instances, the separator is a point
 - In case of 2 dimensional instances, the separator is a line
 - In case of 3 dimensional instances, the separator is a plane
 - In case of instances in more than 3 dimensional space, the separator is a hyperplane
- Given a set of linearly separable instances, there exist infinite number of linear separators which can separate the instances into 2 classes
 - SVM chooses that linear separator which has the maximum "margin"
 - Intuition: A linear separator with the maximum margin will generalize better for new unseen instances in test data

SVM: Linear Separator



- An example of two dimensional instances
- Two classes:
 - Positive and Negative
- Linearly separable data
- Infinite number of linear separators are possible

SVM: Linear Separator



- Goal: To find the linear separator which provides the maximum margin of separation between two classes
- Support vectors: Instances which lie on the margin boundaries

- Any other possible way to find hyperplanes?
- What are the issues?

Representation of a hyperplane

The general form of a hyperplane in an n-dimensional space is given by: $w_1x_1 + w_2x_2 + \cdots + w_nx_n + b = 0$

• Can be expressed in a vector form as: $w^T x + b = 0$ or $w \cdot x + b = 0$



- The vector \boldsymbol{w} is perpendicular to the hyperplane and \boldsymbol{b} is the bias term (controls the offset of the hyperplane from the origin). $w \perp hyperplane$
- $w \cdot x + b > 0$ ightarrow Points on one side of the hyperplane.
 - $w\cdot x+b < 0$ ightarrow Points on the other side.
 - If b=0, the hyperplane passes through the origin.
 - If $b \neq 0$, the hyperplane is shifted away from the origin.

Representation of a hyperplane

The general form of a hyperplane in an n-dimensional space is given $w_1x_1+w_2x_2+\cdots+w_nx_n+b=0$ by:

Can be expressed in a vector form as: $w^T x + b = 0$ or $w \cdot x + b = 0$



- For a particular linear separator, any point x^i lying on the "positive" side will have positive value of

$$w\cdot x^i+b$$

- Also, any point Villa lying on the "negative" side will have negative value of $w\cdot x^i+b$
- E.g., the point (2,1.5) is on positive side of Line 2 (0.5) and on negative side of Line 1 (-0.5)

SVM: Training



Training instances: $\{\langle \mathbf{x}^1, y^1 \rangle, \langle \mathbf{x}^2, y^2 \rangle, \dots \langle \mathbf{x}^N, y^N \rangle\}$

- xⁱ is a point in n-dimensional space
- $y^i \in \{+1, -1\}$ its corresponding true class label
- Goal: To find optimal linear separator which maximizes the margin
- All positive class points (+1 label) must be on or beyond the +1 margin line.
- All negative class points (-1 label) must be on or beyond the -1 margin line.
- Since maximizing the margin is the core objective, adding an arbitrary scaling factor k is unnecessary because it cancels out in optimization.

• For +1 points:
$$w \cdot x^i + b \ge 1$$

• For -1 points:
$$w \cdot x^i + b <= -1$$

Computing Margin Width



Consider two points \mathbf{x}^1 and \mathbf{x}^2 such that they lie on the opposite margins and the vector $\mathbf{x}^2 - \mathbf{x}^1$ is perpendicular to the linear separator The vector \mathbf{w} is also perpendicular to the linear separator margin = d = $w^T(x^2 - x^1)$

$$u = \frac{w(x - x)}{\|w\|}$$

Therefore, by definition,

$$(\mathbf{x}^2 - \mathbf{x}^1) = \lambda \cdot \mathbf{w}$$

 $w^T x^1 + b = -1$

$$w^T x^2 + b = 1$$
$$w^T \cdot (\mathbf{x}^2 - \mathbf{x}^1) = 2$$

Computing Margin Width



| • | Substituting $(\mathbf{x}^2 - \mathbf{x}^1) = \lambda \cdot \mathbf{w}$ |
|--------------------|---|
| | $\mathbf{w}^T \cdot (\mathbf{x}^2 - \mathbf{x}^1) = 2$ |
| | $\lambda \mathbf{w}^T \cdot \mathbf{w} = 2 \Rightarrow \lambda = \frac{2}{\mathbf{w}^T \cdot \mathbf{w}}$ |
| • | Margin width: |
| $\ \mathbf{x}^2\ $ | $-\mathbf{x}^{1} \ = \sqrt{(\mathbf{x}^{2} - \mathbf{x}^{1})^{T} \cdot (\mathbf{x}^{2} - \mathbf{x}^{1})}$ |
| $\ \mathbf{x}^2\ $ | $-\mathbf{x}^{1}\ ^{2} = \lambda^{2}(\mathbf{w}^{T} \cdot \mathbf{w})$ |
| $\ \mathbf{x}^2\ $ | $-\mathbf{x}^{1}\ ^{2} = \frac{4}{(\mathbf{w}^{T} \cdot \mathbf{w})^{2}} (\mathbf{w}^{T} \cdot \mathbf{w})$ |
| $\ \mathbf{x}^2\ $ | $\ -\mathbf{x}^{1}\ ^{2} = \frac{4}{\mathbf{w}^{T} \cdot \mathbf{w}} \propto \frac{1}{\mathbf{w}^{T} \cdot \mathbf{w}}$ |
| | |

SVM: Optimization Problem



- Objective:
 - Maximize the margin

$$\min_{w,b}\left(rac{1}{2}w^Tw
ight)$$

- Subject to the following constraints:
 - Every training instance should lie on the appropriate (positive / negative) side of the linear separator

 $y^i(w^Tx^i+b)\geq 1, \quad orall 1\leq i\leq N$

 $-y^i(w^Tx^i+b)+1\leq 0, \quad orall 1\leq i\leq N$

SVM: Soft-margin Formulation



Objective:

?

 Maximize the margin and minimize the training error

$$\min_{w,b,\xi}\left(rac{1}{2}w^Tw
ight)+C\sum_{i=1}^N\xi_i$$

- Subject to the following constraints:
 - Introducing slack variables so that the constraint is satisfied for training instances lying on incorrect side

$$egin{aligned} y^i(w^Tx^i+b) &\geq 1-\xi_i, \quad orall 1 \leq i \leq N \ -eta_i &\leq 0 \quad orall 1 \leq i \leq N \end{aligned}$$

- High C => Low ξ_i
 Low C => High ξ_i -> This can allow for large outliers
- What happens if we remove

Optimization using Lagrange Multipliers

• One Lagrange multiplier is associated with each distinct constraint

$$L(\alpha_{1}, \cdots, \alpha_{N}, \mu_{1}, \cdots, \mu_{N})$$

$$= \min_{\mathbf{w}, w_{0}, \xi} \left(\frac{1}{2} \mathbf{w}^{T} \cdot \mathbf{w}\right) + C \cdot \sum_{i=1}^{N} \xi_{i}$$

$$+ \sum_{i=1}^{N} \alpha_{i} \cdot \left(-y^{i} \left(\mathbf{w}^{T} \cdot \mathbf{x}^{i} + w_{0}\right) + 1 - \xi_{i}\right)$$

$$+ \sum_{i=1}^{N} \mu_{i} \cdot \left(-\xi_{i}\right)$$
s.t. $\alpha_{i} \ge 0, \mu_{i} \ge 0, \forall_{1 \le i \le N}$

Optimization using Lagrange Multipliers

Differentiating w.r.t.

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^{N} \alpha_i y^i \mathbf{x}^i = 0 \qquad \Rightarrow \mathbf{w}^* =$$

ξi

$$\Rightarrow \mathbf{w}^* = \sum_{i=1}^N \alpha_i y^i \mathbf{x}^i$$

• Differentiating w.r.t. w_0

$$\frac{\partial L}{\partial w_0} = \sum_{i=1}^N -\alpha_i y^i = 0 \implies \sum_{i=1}^N \alpha_i y^i = 0$$

Differentiating w.r.t.

$$\frac{\partial L}{\partial \xi_i} = C - \alpha_i - \mu_i = 0 \implies \mu_i + \alpha_i = C$$



Finally,
$$L(\alpha_1, \cdots, \alpha_N) = \frac{-1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y^i y^j \left(\mathbf{x}^{i^T} \mathbf{x}^j \right) + \sum_{i=1}^N \alpha_i$$

- Dual optimization problem:
 - Objective function:

$$\max_{\alpha_1,\cdots,\alpha_N} \frac{-1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y^i y^j (\mathbf{x}^{i^T} \mathbf{x}^j) + \sum_{i=1}^N \alpha_i$$

M

Subject to the following constraints:

$$\alpha_{i} \geq 0, \mu_{i} \geq 0, \mu_{i} + \alpha_{i} = C, \forall_{i} \text{ and } \sum_{i=1}^{N} \alpha_{i} y^{i} = 0$$

$$\Rightarrow \alpha_{i} \geq 0, \alpha_{i} \leq C, \forall_{i} \text{ and } \sum_{i=1}^{N} \alpha_{i} y^{i} = 0$$

Any Quadratic Programming Solver can be used for solving this

• Vary C and show its impact on decision boundary.

Using SVM for Predictions

- How to predict the class label for a new instance x given a trained SVM
- Primal Form: $\mathbf{w}^T \cdot \mathbf{x} + w_0$
 - Compute $\mathbf{x} = \begin{bmatrix} 2, 4 \end{bmatrix}$; Positive value indicates the positive class and vice versa
 - E.g., $\mathbf{w} = [2, -1]$ and $w_0 = -2$ be the learned parameters
 - For the new instance $\mathbf{w}^T \cdot \mathbf{x} + w_0 = -2$ and hence Negative class is predicted
 - **Dual Form:** • Compute $\mathbf{w}^T \cdot \mathbf{x} + w_0 = \left(\sum_{i=1}^N \alpha_i y^i \mathbf{x}^i\right)^T \mathbf{x} = \sum_{i=1}^N \alpha_i y^i \mathbf{x}^{i^T} \mathbf{x}$
 - Positive value indicates the positive class and vice versa
 - Practically, most of the α_i values are zeros; non-zero only for support vectors



- SVM works well with linearly separable data.
- But, many real-world data are non-linearly separable in nature
- The kernel function implicitly maps data into a higher-dimensional space where it becomes linearly separable.
- Instead of explicitly transforming data into a higher-dimensional space, we use a kernel function to compute the dot product in that space efficiently.

Explicit Transformation: High Computational Cost

If we explicitly map data from a lower-dimensional space R^d to a higher-dimensional space R^D using a feature transformation $\phi(x)$, we need to compute the dot product in the new space.

Given n training samples, each with d features, transforming them explicitly to a higher dimension D takes:

O(nD)

Then, computing the dot product between two transformed vectors $\phi(x^i)$ and $\phi(x^j)$ in R^D requires:

0**(**D**)**

Training SVM typically involves solving a quadratic programming problem, which has a complexity of:

 $O(n^2D)$ (in the worst case)

due to the need to compute pairwise dot products in the high-dimensional space. If *D* is very large (e.g., infinite in the case of the Radial Basis Function (RBF) kernel), the computation becomes impractical.

Kernel Function is the Saviour: Avoids Curse of Dimensionality, Efficient Computation

The kernel function $K(x^i, x^j) = \langle \phi(x^i), \phi(x^j) \rangle$ computes the dot product in high dimensional space in: O(d) because it only depends on the original d-dimensional data.

The SVM training complexity **remains**: O(n²d)

instead of $O(n^2D)$, making it computationally feasible even when D is very large or infinite.

| Kernel Type | Equation | Description |
|--|---|--|
| Linear Kernel | $K(x_i,x_j)=x_i\cdot x_j$ | Used when data is already linearly separable. |
| Polynomial Kernel | $egin{aligned} K(x_i,x_j) &= (x_i \cdot x_j + c)^d \end{aligned}$ | Maps data into a higher-degree polynomial space. |
| Radial Basis Function (RBF) Kernel | $egin{aligned} K(x_i,x_j) = \ \exp(-\gamma \ x_i-x_j\ ^2) \end{aligned}$ | Maps data to an infinite-dimensional space. Useful for complex boundaries. |
| Sigmoid Kernel | $egin{aligned} K(x_i,x_j) = \ 	anh(lpha x_i \cdot x_j + c) \end{aligned}$ | Similar to neural network activation functions. |

Linear Kernel

$$K(x_i,x_j)=x_i\cdot x_j$$

Used when data is already linearly separable.

 $K(x^{i},x^{j})=x_{i}^{T}x_{j}$

- Directly computes the dot product in O(d).
- No need for explicit mapping.

| Polynomial Kernel | $K(x_i,x_j)=(x_i \ \cdot$ | Maps data into a higher-degree |
|-------------------|---------------------------|--------------------------------|
| | $(x_j+c)^d$ | polynomial space. |

 $\mathsf{K}(\mathsf{x}_{i'},\mathsf{x}_{j}) = (\mathsf{x}_{i}^{\mathsf{T}}\mathsf{x}_{j} + \mathsf{C})^{\mathsf{p}}$

- Requires computing x_i⁻x_j (which takes O(d)) and then raising it to power p (which takes O(1)).
- Total complexity: **O(d)**.

How the Kernel Trick Works in SVM?

- 1. Choose an appropriate kernel function based on the dataset.
- 2. Compute the kernel matrix $K(x_i, x_j)$, which represents inner products in higher-dimensional space.
- 3. Use this matrix in the SVM optimization problem without explicitly transforming the data.
- 4. The SVM classifier finds the optimal hyperplane in the transformed space.

Advantage of Dual over Primal

Handles High-Dimensional Data Efficiently (Kernel Trick)

- The **primal form** works directly in the original feature space, making it impractical when the number of features (D) is large or infinite (e.g., in kernelized SVM).
- The **dual form** allows the use of the **kernel trick**, enabling SVMs to operate in an **implicit high-dimensional space** without explicitly transforming data.

$$L(\alpha_{1}, \dots, \alpha_{N}, \mu_{1}, \dots, \mu_{N})$$

$$= \min_{\mathbf{w}, w_{0}, \xi} \left(\frac{1}{2} \mathbf{w}^{T} \cdot \mathbf{w}\right) + C \cdot \sum_{i=1}^{N} \xi_{i}$$

$$+ \sum_{i=1}^{N} \alpha_{i} \cdot \left(-y^{i} (\mathbf{w}^{T} \cdot \mathbf{x}^{i} + w_{0}) + 1 - \xi_{i}\right)$$

$$+ \sum_{i=1}^{N} \mu_{i} \cdot \left(-\xi_{i}\right)$$

$$\alpha_{1}, \dots, \alpha_{N} \frac{-1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_{i} \alpha_{j} y^{i} y^{j} (\mathbf{x}^{i^{T}} \mathbf{x}^{j}) + \sum_{i=1}^{N} \alpha_{i}$$

Primal

Dual