

Kernel machines

Kernel machines are a class of algorithms that leverage kernel functions to create non-linear decision boundaries allowing for complex data analysis.

The main use of the kernel is to transforming data into high dimensions,

They are mostly effective in both classification as well as regression

Definitions of Kernel functions

→ A Kernel function $K(x, y)$, calculates the inner product of 2 points x and y in a high-dimensional feature space

common forms

- > linear kernel
- > polynomial Kernel
- > Radial Basis Functions (RBF)
- > sigmoid Kernel

\rightarrow sigmoid ...

Properties of Kernel

- \rightarrow symmetry
- \rightarrow positive semidefiniteness
- \rightarrow Inner product representations

* Roles of Kernel in non linear data

Kernel transform non-linear separable data into higher-dimensional space.

- \rightarrow linear Kernel - $K(x, y) = x \cdot y$
- \rightarrow polynomial $K(x, y) = (x \cdot y + c)^d$
- \rightarrow RBF $K(x, y) = e^{-\gamma \|x - y\|^2}$
- \rightarrow Sigmoid Kernel $K(x, y) = \tanh(\alpha \cdot x \cdot y + c)$

Intro to Kernels for Structured data & text

- \rightarrow Kernel function compute similarities between data point in high-dimensional space

' Point in high-dimensional space

- Handling structured and unstructured data
- NLP, Web documents, social networks.

Structured Data

- Data with predefined schema
- tree, graphs
- cons:- finding relationship

~~Text~~

- unstructured data, represented as sequence of word or characters
- Articles, emails, and social media posts
- cons:- high dimensionality and sparsity.

Examples

- ① Document classification
- ② sentiment analysis
- ③ spam detection

Q Biological sequence

why structured Data need special Kernels?

- traditional Kernel fails to capture hierarchical relation ship structure

why text Data need special Kernels?

→ High Dimensionality

→ sparsity : most word or tokens in corpora are infrequent

→ semantic meaning :- words with similar meaning may not share the same representation

→ order sensitivity :- word order impact meaning, requires kernels that account for sequence

→ Heterogeneous Data : text data can vary significantly in length and structure.

Story Kernels :- measure the similarity based on common substory

TF-IDF Kernel :- Based on term frequency-inverse

Word Embedding Kernel : Use embedding like word 2 vec

* by one for similarities

Multiple Kernel learning

- > combines multiple kernel functions to improve predictions
- > Help integrate different data type
- > used in bioinformatics, computer vision, and NLP

Why multiple kernels

- > A single kernel may not capture data complexity
- > Different datatypes need different kernels
- > improves accuracy

core concepts :-

Kernel functions:- Measure similarity between data points.

linear combinations : uses weighted sum of multiple kernels

Optimizations: learns the best kernel mix
and classifier parameters.

Input: taking multiple datatypes
apply different kernel function
combine kernels using weights
optimize weights and classified
final predictive model.

Linear Kernel :- capture direct relationship

RBF Kernel : handles nonlinear data

Polynomial Kernel : models complex patterns

Pros :-

- capture true data feature / types
- > improves prediction accuracy
- > works with multiple data sources
- high computational cost
- complex optimizations

- Complex optimizations
- Risk of overfitting w/ too many kernels

Emerging trends in MKL

- Sparse MKL: uses only the most relevant kernels
- Hierarchical MKL: layered kernel approach
for complex data
- Deep Kernel learning: combines MKL with deep learning

Applications

- Bioinformatics
- computer vision
- NLP
- Health care analytics

Case Study Exam

Sentiment analysis w/ customer rev.

RBF Kernel: used for numerical data

String Kernel: used for text review

Model combination: improve accuracy by
stacking data sets

MOD2

02 February 2025 11:20

A real-time system processes a data stream. The first three observations are as follows:

Point	x_1	x_2	Label (y)
A	2	3	+1
B	-1	-2	-1
C	0	1	+1

$$n = 0.4 \quad \lambda = 0.1$$

$$y \cdot (w \cdot x + b) \geq 1$$

update weight

$$\leftarrow w + n(y \cdot x - \lambda \cdot w)$$

$$(2, 3) \quad y = +1$$

$$= 1([0, 0] \cdot [2, 3] + 0) = 0$$

$$= 0$$

$$x = y, x \rightarrow w$$

$$a = 1 \cdot (2, 3) - 0.1 [0, 0]$$

$$= (2, 3)$$

$$w + n(a)$$

$$= (0, 0) + 0.4 [2, 3]$$

$$w = (0.2, 0.3)$$

-- $w_i, b \rightarrow$

$$b = b + \eta y$$

$$\begin{aligned} b &= 0 + 0.1 / \\ &= 0.1 \end{aligned}$$

for $i = 2$

$$x_i = [-1, -2] \quad y_i = -1$$

$$y_i(w \cdot x + b)$$

$$\rightarrow (0.2, 0.3)([-1, -2]) + 0.1$$

$$= -1(0.9 + 0.1)$$

$$= -1(-0.5)$$

$$= 0.75$$

$$w = w + \eta(y_i x_i - \lambda w)$$

$$a = y_i x_i - \lambda w$$

$$= -1[-1, -2] - 0.1[0.2, 0.3]$$

$$= (1, 2) - (0.01, 0.03)$$

$$a = (0.99, 1.97)$$

$$w = w + \eta(a)$$

$$= [0.2, 0.3] + 0.1[0.99, 1.97]$$

$$= [0.2, 0.3] + 0.09, 0.19]$$

$$w = [0.29, 0.47]$$

$$b = b + \eta y$$

$$b = 0.1 + 0.1(-1)$$

$$= 0.1 - 0.1$$

$$b = 0$$

$$x_3 = [0, 1] \quad y_3 = +1)$$

$$y \cdot (\omega x + b) < 1$$

$$1([0.3, 0.5][0, 1] + 0) >$$

$$0.5 < 1$$

$$\omega = \omega + n(\alpha)$$

$$\alpha = y_n - \gamma \omega$$

$$= 1(0, 1) - 0.1(0.3, 0.5)$$

$$= (0, 1) - (0.03, 0.05)$$

$$= (0, 0.95)$$

$$\omega = (0.3, 0.5) + 0.1(0, 0.95)$$

$$\boxed{\omega = (0.3, 0.595)}$$

$$y \cdot (\omega x + b) > 1$$

$$1([0.3, 0.6][2, 3] + 0) > 1$$

final weights $[0.3, 0.6]$

$$0.3x_1 + 0.6x_2 = 0$$

$$\boxed{x_1 + 2x_2 = 0} \text{ Hyperplane}$$

$$\text{margin } \frac{2}{\|\omega\|}$$

$$\|\omega\| = \sqrt{0.09 + 0.36}$$

$$= \sqrt{0.49}$$

$$\approx 0.67$$

$$\approx \frac{2}{0.67} = 2.98$$

Theory

Hard margin SVM

- to find a hyperplane that perfectly separate the classes in dataset
- used when data is linearly separable with out overlap

Types

- one type zero misclassification

Working

- The SVM tries to find hyperplane that maximize the margin
- Since misclassification is not allowed data points must classified correctly
- It uses some constraints that no point can be within the margin or misclassified.

Advantages

→ works well the data is perfectly
separable

→ maximizes the margin

Disadvantage

→ fails when data has noise or overlap

→ sensitive and slow and
misclassifies

Soft margin SVM

Goal

→ To allow some misclassification in order to
create a more generalizable model

→ used when data is not perfectly separable

Types

→ linear soft margin: Allow some misclassification
while maintaining a linear decision boundary

→ non-linear soft margin SVM: Use kernel
to handle complex data distributions

Working

→ Introduces slack variables (ξ) to allow some
misclassified points

- Balance between maximizing margin and minimizing classification error
- uses a regularization parameter C to control the the trade off between a wider margin & few misclassifications

Advantages :-

- Handles noisy and overlapping data better
- more flexible in real world applications
- can work with non linear data using Kernel

Disadvantage :-

- needs careful C parameter
- computational expensive

Feature	Hard Margin SVM	Soft Margin SVM
Misclassification	Not allowed (strict)	Allowed (with penalties)
Data Type	Works only for perfectly separable data	Works for overlapping or noisy data
Slack Variables (ξ)	Not used	Used to allow margin violations
Outlier Sensitivity	Highly sensitive	More robust
Generalization	Poor for real-world data	Better generalization
Computational Cost	Lower	Higher (due to slack variables and optimization)

Distributed SVM

→ handles the date or large scale data
 that can't fit into memory of single
 machine . It used when data has
 to spread across diff locations

Challenges

- Data heterogeneity
- & communication overhead
- Scalability
- Data imbalance

Algorithms

- Parallel SVM
- MapReduce SVM
- Gradient Based Distribution

work flow

- Data partitioning : Divide the data across the multiple machines
- Parallel training
- Model Aggregation
- Model evaluation
- Model update

Hyperparameter tuning

- Regularization parameter C :
- Ctrl trade off b/w min & max misclassification
- Kernel parameters :
- learning rate
- communication setting : tune communication intervals b/w nodes. To minimize overhead

Read 1-10.

Applications

- Large Scale Image Classification
- Big Data Analytics
- Text mining
- Bioinformatics
- Autonomous System

Online SVM

→ Handles the streaming data, unlike regular SVM, where entire dataset has to be stored in memory for training. Online SVM update the model incrementally as new data points arrive, making them suitable for situations.

Challenges

- Memory usage
- Model drift
- Computational complexity
- Model convergence
- Hyperparameter tuning

Workflow

- Data streaming
- Incremental training
- Model update
- Model evaluation
- Drift detection

tuning

- Reg Parameter C
- learning rate η
- kernel parameters
- update frequency

Applications of online SVM

- Real-time classifiers
- Predictive maintenance
- Spam filters
- social media analysis
- sensor networks

Gradient Boosting / AdaBoost