

Kernel methods

Narayana Santhanam

EE 645

Jan 22, 2023

Linear \rightarrow Non-linear

Two approaches:

Kernel methods

Linear \rightarrow Non-linear

Two approaches:

Kernel methods

guarantees, well understood

Linear → Non-linear

Two approaches:

Kernel methods

- guarantees, well understood
- computationally intensive (small data)

Neural networks

- less well understood

Linear → Non-linear

Two approaches:

Kernel methods

- guarantees, well understood
- computationally intensive (small data)

Neural networks

- less well understood
- computationally cheap (large data)

Next couple of weeks

Kernel methods

Next couple of weeks

Kernel methods

High level picture

Next couple of weeks

Kernel methods

High level picture

Support Vector classification and regression

Next couple of weeks

Kernel methods

High level picture

Support Vector classification and regression

Kernel PCA

Next couple of weeks

Kernel methods

High level picture

Support Vector classification and regression

Kernel PCA

Random Fourier Futures

Next couple of weeks

Kernel methods

High level picture

Support Vector classification and regression

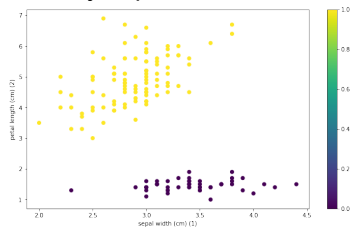
Kernel PCA

Random Fourier Features

Credit risk, Power data

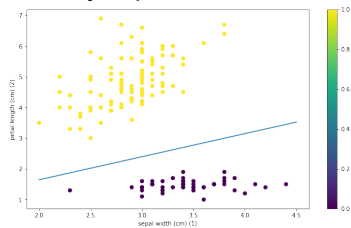
Classification

Linearly separable



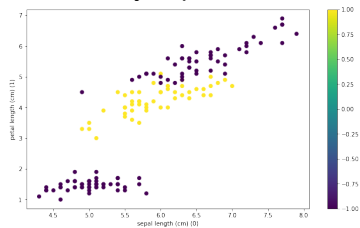
Classification

Linearly separable



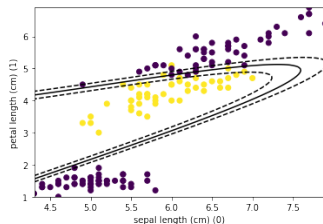
Classification

Not linearly separable



Classification

Not linearly separable



Boundaries not linear..

but are linear in a higher dimensional space!

Closer look

Principled approach for linear \rightarrow nonlinear

Powerful, yet generalizes well

Often explainable

at least more than other state of art

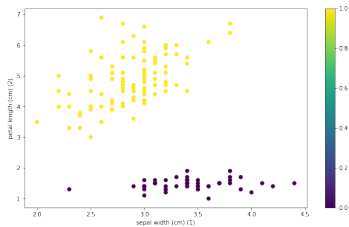
New advances increase reach (more data)

but not as much as NN

Stepping back

Linearly separable

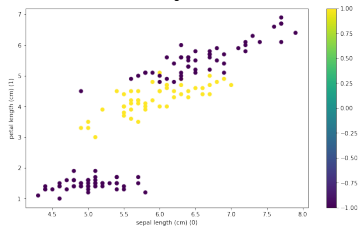
Where would you draw a *linear* classifier?



Stepping back

Not linearly separable

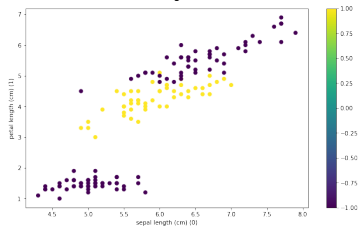
Where would you draw a *linear* classifier?



Stepping back

Not linearly separable

Where would you draw a *linear* classifier?

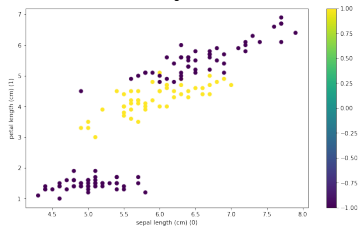


Minimum number of mistakes?

Stepping back

Not linearly separable

Where would you draw a *linear* classifier?

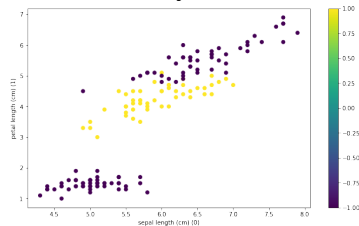


Minimum number of mistakes?
infeasible, NP-hard

Stepping back

Not linearly separable

Where would you draw a *linear* classifier?



Minimum number of mistakes?

infeasible, NP-hard

Instead: variation of ℓ_1 loss, sum of margins

Setting up linear classification

Distance of x from a plane $w^T x - b = 0$?

Setting up linear classification

Distance of x from a plane $w^T x - b = 0$?

$$\frac{w^T x - b}{\|w\|}$$

Formulation of Support Vector Classification

Analyze this in some detail

Formulation of Support Vector Classification

Analyze this in some detail

One of the key advantages: interpretability

Comes through formulation

Formulation of Support Vector Classification

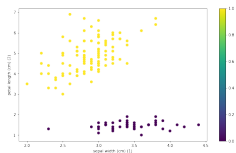
Analyze this in some detail

One of the key advantages: interpretability

Comes through formulation

Primal/Dual formulation is a key optimization idea
accelerations of training neural networks
resource allocation/optimization in economics,
urban planning

Formulation



Linearly separable points:
Training data:

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

Margin (closest distance from separating boundary)

$$\gamma(w, b) = \min_{x_i} \frac{|w^T x_i - b|}{\|w\|}$$

Redundant parameterization scaling w, b by any number does not change the margin

Formulation

Support vector machine formulation:

$$w^*, b^* = \arg \max_{w, b} \gamma(w, b)$$

subject to $y_i(w^T x_i - b) \geq 0$ (for all training (x_i, y_i))

Redundant parameterization: scaling w, b changes nothing

- only the directions/intercepts matter
- represent by scaling of w, b that ensures

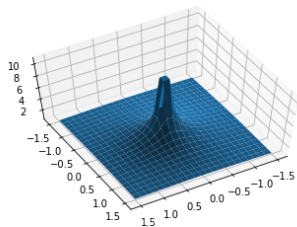
$$\min_{x_i} |w^T x_i - b| = 1$$

Formulation

Support vector machine formulation
(removing the redundant formulation)

$$w^*, b^* = \arg \max_{w, b} \frac{1}{\|w\|}$$

subject to $y_i(w^T x_i - b) \geq 1$ (for all training (x_i, y_i))



Formulation

$\frac{1}{\|w\|}$ not convex/concave

But it is easy to come with a convex formulation

$$w^*, b^* = \arg \min_{w, b} \frac{1}{2} \|w\|^2$$

subject to $y_i(w^T x_i - b) \geq 1$ (for all pairs (x_i, y_i))

A bit on convex optimization

General concepts beyond support vector machine formulation

A bit on convex optimization

General concepts beyond support vector machine formulation

Lagrangian:

$$L(w, b, \Lambda) = \frac{1}{2} \|w\|^2 + \sum_i \lambda_i (1 - y_i(w^T x_i - b))$$

$L(w, b, \Lambda)$ is a key idea in any constrained optimization

Primal/dual formulation

Neat property of Lagrangians:

Optimal point

$$\min_{w,b} \max_{\Lambda \geq 0} L(w, b, \Lambda) \text{ (primal formulation)}$$

Dual formulation (Nash, von Neumann)

$$\max_{\Lambda \geq 0} \min_{w,b} L(w, b, \Lambda) \text{ (dual formulation)}$$

Primal/dual formulation

Neat property of Lagrangians:

Optimal point

$$\min_{w,b} \max_{\Lambda \geq 0} L(w, b, \Lambda) \text{ (primal formulation)}$$

Dual formulation (Nash, von Neumann)

$$\max_{\Lambda \geq 0} \min_{w,b} L(w, b, \Lambda) \text{ (dual formulation)}$$

Incidentally Nash won both the Nobel Prize in Econ and the Abel Prize

Primal/dual formulation

For free:

$$\min_{w,b} \max_{\Lambda \geq 0} L(w, b, \Lambda) \geq \max_{\Lambda \geq 0} \min_{w,b} L(w, b, \Lambda)$$

But equality in many cases

- in many convex formulations, including our current case
- so one could solve either version
- dual in kernel methods very insightful for explainability

Dual formulation: 2 key insights

- Setting gradient of $L(w, b, \Lambda)$ to 0, optimal w^* satisfies

$$w^* = \sum_i \lambda_i y_i x_i$$

Representer theorem: solution w^* is linear combination of inputs

Dual formulation: 2 key insights

- Setting gradient of $L(w, b, \Lambda)$ to 0, optimal w^* satisfies

$$w^* = \sum_i \lambda_i y_i x_i$$

Representer theorem: solution w^* is linear combination of inputs

- Finding Λ s

$$\max_{\Lambda \geq 0} \sum \lambda_i - \frac{1}{2} [\lambda_1 y_1 \quad \dots \quad \lambda_n y_n] X X^T \begin{bmatrix} \lambda_1 y_1 \\ \vdots \\ \lambda_n y_n \end{bmatrix}$$

Data only shows up through dot products (XX^T)
Crux of Kernel approach to nonlinearity