# Foundations of Machine Learning AI2000 and AI5000

FoML-35
Support Vector Machines (cntd.)
Duality to obtain the max margin classification

Dr. Konda Reddy Mopuri
Department of AI, IIT Hyderabad
July-Nov 2025

భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

DiL
Data-driven Intelligence
& Learning Lab

# So far in FoML

- Intro to ML and Probability refresher

- MLE, MAP, and fully Bayesian treatment

- Supervised learning

  a. Linear Regression with basis functions

  b. Bias-Variance Decomposition

  c. Decision Theory - three broad classification strategies

  d. Neural Networks

- Unsupervised learning

  a. K-Means, Hierarchical, and GMM for clustering

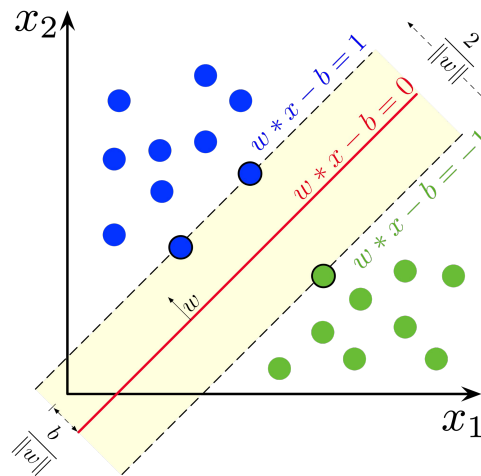- Kernelizing linear Models

  a. Dual representation, Kernel trick

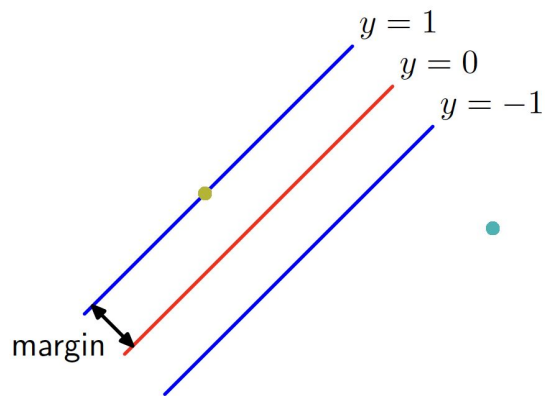భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
**Indian Institute of Technology Hyderabad**

DiL
Data-driven Intelligence
& Learning Lab

# For today

- SVM (cntd.)
  - Duality to obtain the max margin classification

భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

DiL
Data-driven Intelligence
& Learning Lab

# Max margin classifier

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

$$t_n\left(\mathbf{w}^{\mathrm{T}}\phi(\mathbf{x}_n) + b\right) \geqslant 1,$$

$$n = 1, \ldots, N.$$

# Max margin classifier

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \qquad t_n \left(\mathbf{w}^{\mathrm{T}} \phi(\mathbf{x}_n) + b\right) \geqslant 1, \qquad n = 1, \ldots, N.$$

$$\left[ t_n y(\mathbf{x}_n) - 1 \right] \geqslant 0$$

$$f(x) \qquad - a_n \, g(x)$$

- Primal Lagrangian

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^{N} a_n \left\{ t_n (\mathbf{w}^{\mathrm{T}} \phi(\mathbf{x}_n) + b) - 1 \right\}$$

$$L(w, \lambda) \qquad f(w) \sim \qquad g(w)$$
$$b$$

భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

DiL
Data-driven Intelligence
& Learning Lab

# Max margin classifier

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2}\|\mathbf{w}\|^2 - \sum_{n=1}^{N} a_n \left\{ t_n(\mathbf{w}^{\mathrm{T}}\phi(\mathbf{x}_n) + b) - 1 \right\}$$

KKT conditions

$$t_n(\mathbf{w}^T\mathbf{x}_n + b) - 1 \geq 0 \qquad \text{for} \qquad n = 1,\dots,N$$

$$a_n \geq 0 \qquad \text{for} \qquad n = 1,\dots,N$$

$$a_n(t_n(\mathbf{w}^T\mathbf{x}_n + b) - 1) = 0 \qquad \text{for} \qquad n = 1,\dots,N$$

# Max margin classifier

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2}\|\mathbf{w}\|^2 - \sum_{n=1}^{N} a_n \left\{ t_n(\mathbf{w}^\mathrm{T}\phi(\mathbf{x}_n) + b) - 1 \right\}$$

KKT conditions

$$t_n(\mathbf{w}^T\mathbf{x}_n + b) - 1 \geq 0 \qquad \text{for} \qquad n = 1, \ldots, N$$

$$a_n \geq 0 \qquad \text{for} \qquad n = 1, \ldots, N$$

$$a_n(t_n(\mathbf{w}^T\mathbf{x}_n + b) - 1) = 0 \qquad \text{for} \qquad n = 1, \ldots, N$$

Derive the dual Lagrangian via

$$\frac{\partial L}{\partial \mathbf{w}} = 0, \quad \frac{\partial L}{\partial b} = 0 \implies \tilde{L}(\mathbf{a}) = \min_{\mathbf{x}, b} L(\mathbf{x}, b, \mathbf{a})$$

$w$

$b$

# Max margin classifier

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2}\|\mathbf{w}\|^2 - \sum_{n=1}^{N} a_n \left\{ t_n(\mathbf{w}^{\mathrm{T}}\phi(\mathbf{x}_n) + b) - 1 \right\}$$

KKT conditions

$$t_n(\mathbf{w}^T\mathbf{x}_n + b) - 1 \geq 0 \qquad \text{for} \qquad n = 1,\ldots,N$$

$$a_n \geq 0 \qquad \text{for} \qquad n = 1,\ldots,N$$

$$a_n(t_n(\mathbf{w}^T\mathbf{x}_n + b) - 1) = 0 \qquad \text{for} \qquad n = 1,\ldots,N$$

$\tilde{L}(a)$

Derive the dual Lagrangian via $\quad \dfrac{\partial L}{\partial \mathbf{w}} = 0, \quad \dfrac{\partial L}{\partial b} = 0 \quad \Longrightarrow \quad \tilde{L}(\mathbf{a}) = \min_{\mathbf{x},b} L(\mathbf{x}, b, \mathbf{a})$

Now, solve for a* $\quad \mathbf{a}^* = \arg\max_{\mathbf{a}} \tilde{L}(\mathbf{a})$

భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

DiL
Data-driven Intelligence
& Learning Lab

# Max margin classifier

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2}\|\mathbf{w}\|^2 - \sum_{n=1}^{N} a_n \left\{ t_n(\mathbf{w}^\mathrm{T}\phi(\mathbf{x}_n) + b) - 1 \right\}$$

KKT conditions

$$t_n(\mathbf{w}^T\mathbf{x}_n + b) - 1 \geq 0 \qquad \text{for} \quad n = 1,\dots,N$$

$$a_n \geq 0 \qquad \text{for} \quad n = 1,\dots,N$$

$$a_n(t_n(\mathbf{w}^T\mathbf{x}_n + b) - 1) = 0 \quad \text{for} \quad n = 1,\dots,N$$

Derive the dual Lagrangian via $\quad \dfrac{\partial L}{\partial \mathbf{w}} = 0, \quad \dfrac{\partial L}{\partial b} = 0 \quad \Longrightarrow \quad \tilde{L}(\mathbf{a}) = \min_{\mathbf{x},b} L(\mathbf{x}, b, \mathbf{a})$

Now, solve for a*  $\quad \mathbf{a}^* = \arg\max_{\mathbf{a}} \tilde{L}(\mathbf{a}) \quad$ then, solve for w*, b*  $\quad \mathbf{w}^*, b^* = \arg\min_{\mathbf{w},b} L(\mathbf{w}, b, \mathbf{a}^*)$

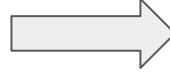# Max margin classifier

- Let's form the dual Lagrangian for $L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2}\|\mathbf{w}\|^2 - \sum_{n=1}^{N} a_n \left\{ t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - 1 \right\}$

బరువు

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w}^T - \sum_{n=1}^{N} a_n t_n \mathbf{x}_n^T = 0 \implies \mathbf{w} = \sum_{n=1}^{N} a_n t_n \mathbf{x}_n$$

$x_n \text{ or } \phi(x_n)$

$$\frac{\partial L}{\partial b} = -\sum_{n=1}^{N} a_n t_n = 0 \implies \sum_{n=1}^{N} a_n t_n = 0$$

Eliminate w and b from L

భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

DiL
Data-driven Intelligence
& Learning Lab

# Max margin classifier

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2}\|\mathbf{w}\|^2 - \sum_{n=1}^{N} a_n \left\{ t_n(\mathbf{w}^{\mathrm{T}}\phi(\mathbf{x}_n) + b) - 1 \right\}$$

Applying the stationarity conditions

$$\mathbf{w} = \sum_{n=1}^{N} a_n t_n \mathbf{x}_n \qquad \sum_{n=1}^{N} a_n t_n = 0$$

$$\tilde{L}(\mathbf{a}) = \frac{1}{2} w^T w - \sum_{n=1}^{N} a_n t_n w^T \phi(x_n) - \sum_{n=1}^{N} a_n t_n b + \sum_{n=1}^{N} a_n$$

$$= w^T \left[ \frac{1}{2} w - \sum_{n=1}^{N} a_n t_n \phi(x_n) \right] - b \cdot \sum_{n=1}^{N} a_n \cdot t_n + \sum_{n=1}^{N} a_n$$

$$= w^T \left( -\frac{1}{2} w \right) + \sum_{n=1}^{N} a_n$$

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^{N} a_n - \frac{1}{2}\sum_{n=1}^{N} \sum_{m=1}^{N} a_n a_m t_n t_m \mathbf{x}_n^T \mathbf{x}_m$$

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^{N} a_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

# Max margin classifier

- Dual representation of the max margin (maximize w.r.t **a**)

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^{N} a_n - \sum_{n=1}^{N} \sum_{m=1}^{N} a_n a_m t_n t_m \mathbf{x}_n^T \mathbf{x}_m$$

Such that

$$a_n \geq 0 \ \forall n = 1, \ldots N \qquad \sum_{n=1}^{N} a_n t_n = 0$$

✓ It's a quadratic optimization problem
linear constraints ⟹ Convex region ⟹ local optima = global

✓ However, because of Complexity, in practice we use decomposition techniques (chunking, smo)

we can apply kernel trick

$$K(\underline{x_n}, \underline{x_m})$$
↓
Can now learn Complex nonlinear decision boundary

# Max margin classifier

- New prediction $y(\mathbf{x}_n) = \mathbf{w}^T \mathbf{x}_n + b$ $\implies$ $y(\mathbf{x}) = \sum_{n=1}^{N} a_n t_n \mathbf{x}_n^T \mathbf{x} + b$

# Max margin classifier

- New prediction $y(\mathbf{x}_n) = \mathbf{w}^T \mathbf{x}_n + b$ $\implies$ $y(\mathbf{x}) = \sum_{n=1}^{N} a_n t_n \mathbf{x}_n^T \mathbf{x} + b$

$$\begin{cases} t_n(\mathbf{w}^T\mathbf{x}_n + b) - 1 \geq 0 & \text{for} \quad n = 1, \ldots, N \\ a_n \geq 0 & \text{for} \quad n = 1, \ldots, N \\ a_n(t_n(\mathbf{w}^T\mathbf{x}_n + b) - 1) = 0 & \text{for} \quad n = 1, \ldots, N \end{cases}$$

భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

DiL
Data-driven Intelligence
& Learning Lab

# Max margin classifier

- New prediction $y(\mathbf{x}_n) = \mathbf{w}^T \mathbf{x}_n + b$ $\Longrightarrow$ $y(\mathbf{x}) = \sum_{n=1}^{N} a_n t_n \mathbf{x}_n^T \mathbf{x} + b$

$$t_n(\mathbf{w}^T \mathbf{x}_n + b) - 1 \geq 0 \qquad \text{for} \qquad n = 1, \ldots, N$$
$$a_n \geq 0 \qquad \text{for} \qquad n = 1, \ldots, N$$
$$a_n(t_n(\mathbf{w}^T \mathbf{x}_n + b) - 1) = 0 \qquad \text{for} \qquad n = 1, \ldots, N$$

- Consider $a_n$
  - $> 0 \rightarrow$ lie at margin distance $\rightarrow$ support vectors
  - $= 0 \leftarrow$ lie far from classifier

$$t_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) = 1$$

 భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

DiL
Data-driven Intelligence
& Learning Lab

# Max margin classifier

- New prediction $y(\mathbf{x}_n) = \mathbf{w}^T \mathbf{x}_n + b$ $\implies$ $y(\mathbf{x}) = \sum_{n=1}^{N} a_n t_n \mathbf{x}_n^T \mathbf{x} + b$

$\implies$ $y(\mathbf{x}) = \sum_{m \in S} a_m t_m k(\mathbf{x}_m, \mathbf{x}) + b$

- Find b using $t_n y_n(\mathbf{x}) = 1$ for support vectors

$$t_n \left( \sum_{m \in S} a_m t_m k(\mathbf{x}_m, \mathbf{x}_n) + b \right) = 1$$

$$\sum_{m \in S} a_m t_m k(\mathbf{x}_m, \mathbf{x}_n) + b = t_n$$

$$b = t_n - \sum_{m \in S} a_m t_m k(x_m, x_n)$$

We can consider the average of multiple such estimates (one for a support vector)

# Next

- Gaussian Processes