

SUPPORT VECTOR MACHINE OPTIMIZATION

Pei-Yuan Wu
Electrical Engineering Department
National Taiwan University



OUT-LINE

- Loss function for SVM
- Gradient Descent Algorithm

OPTIMIZATION

- Learning an SVM has been formulated as a constrained optimization problem over w and ξ

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}, \xi_1, \dots, \xi_N \geq 0} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i$$

subject to $y_i(w^T x_i + b) \geq 1 - \xi_i$, for $i = 1, \dots, N$

- The constraint $y_i(w^T x_i + b) \geq 1 - \xi_i$, can be written more concisely as

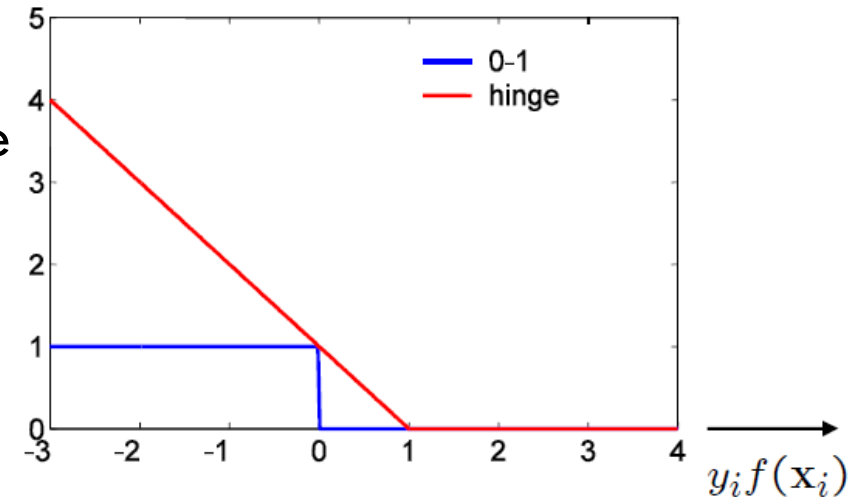
$$y_i f(x_i) \geq 1 - \xi_i$$

which, together with $\xi_i \geq 0$, is equivalent to

$$\xi_i \geq \max(0, 1 - y_i f(x_i))$$

- Hence the learning problem is equivalent to the unconstrained optimization problem over w

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}} \underbrace{\frac{1}{2} \|w\|^2}_{\text{regularization}} + C \underbrace{\sum_{i=1}^N \max(0, 1 - y_i f(x_i))}_{\text{hinge loss}}$$

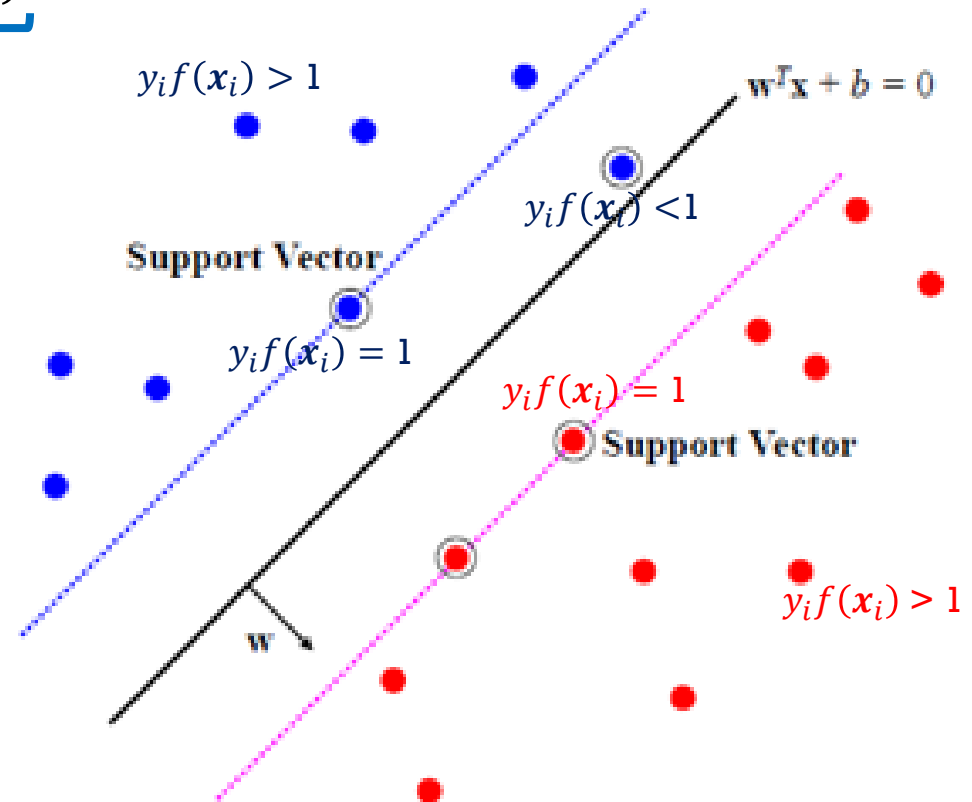


LOSS FUNCTION

$$\min_{w \in \mathbb{R}^d} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \underbrace{\max(0, 1 - y_i f(x_i))}_{\text{hinge loss}}$$

Points are in three categories:

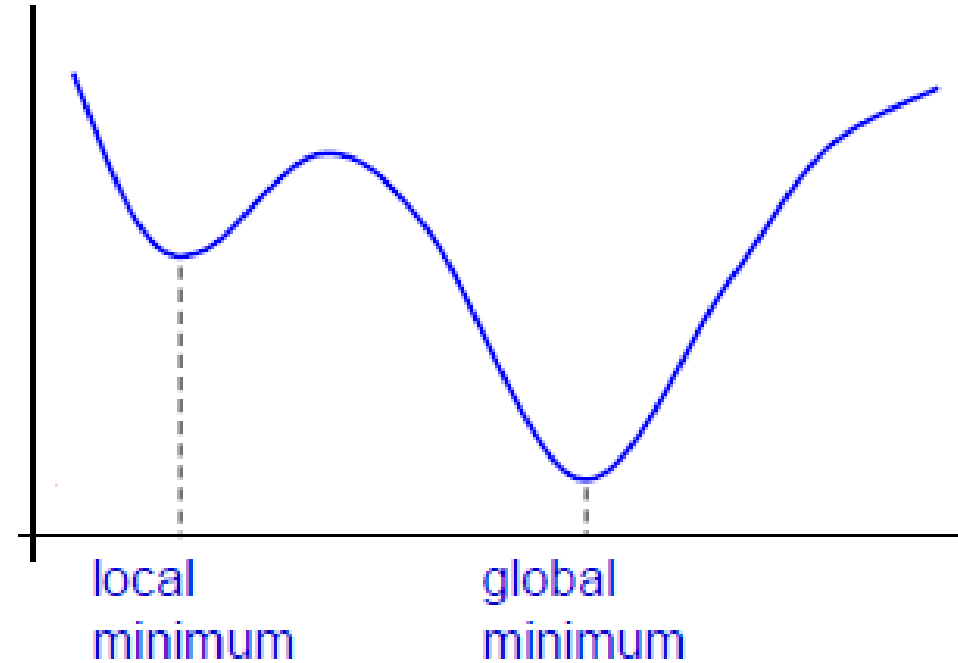
- $y_i f(x_i) > 1$
 - Point is outside margin.
 - No contribution to loss
- $y_i f(x_i) = 1$
 - Point is on margin.
 - No contribution to loss.
- $y_i f(x_i) < 1$
 - Point violates margin constraint.
 - Contributes to loss



OPTIMIZATION CONTINUED

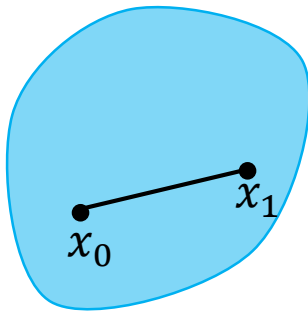
$$\min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \max(0, 1 - y_i f(\mathbf{x}_i))$$

- Does this loss function have a unique minimal solution?
- Does the solution depend on the starting point of an iterative optimization algorithm (such as gradient descent)?

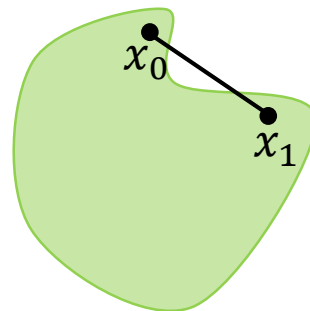


CONVEX FUNCTIONS

We say $\Omega \subset \mathbb{R}^d$ is convex if every $x_0, x_1 \in \Omega$ and $0 \leq t \leq 1$ satisfy $(1-t)x_0 + tx_1 \in \Omega$

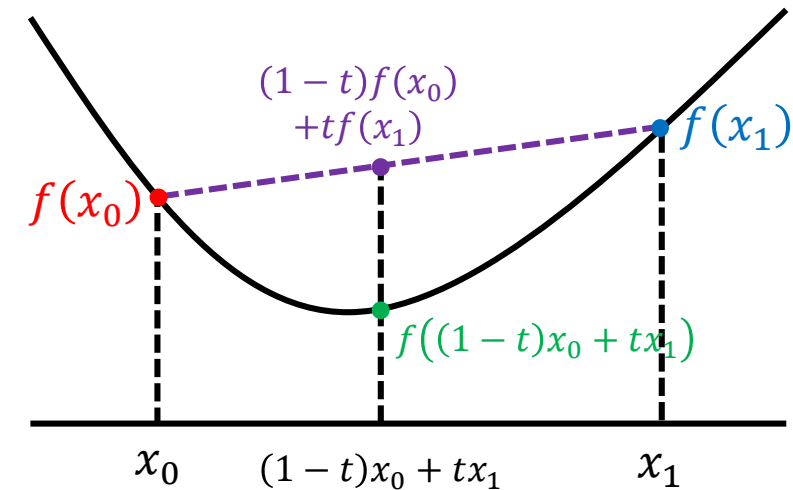


convex



non-convex

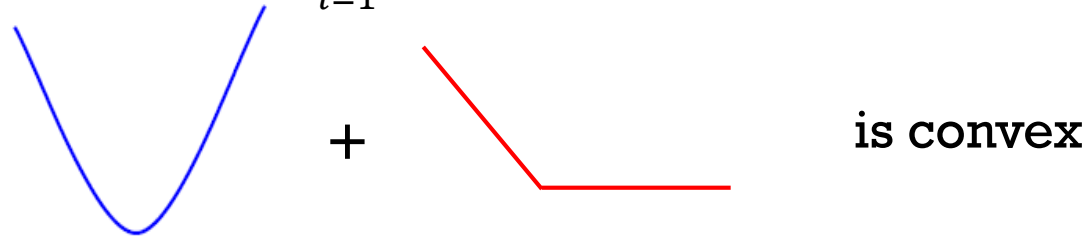
Let Ω be a convex set in \mathbb{R}^d . A function $f: \Omega \rightarrow \mathbb{R}$ is convex if every $x_0, x_1 \in \Omega$ and $0 \leq t \leq 1$ satisfy $f((1-t)x_0 + tx_1) \leq (1-t)f(x_0) + tf(x_1)$



Line joining $(x_0, f(x_0))$ and $(x_1, f(x_1))$ lies above the function graph

CONVEX FUNCTION PROPERTIES

SVM's loss function: $\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b))$



Lemma: Let Ω be a convex set in \mathbb{R}^d , f, g be convex functions on Ω , $\alpha, \beta \geq 0$, then $\alpha f + \beta g$ is convex.

Proof: For every $x_0, x_1 \in \Omega$ and $0 \leq t \leq 1$, one has

$$f((1-t)x_0 + tx_1) \leq (1-t)f(x_0) + tf(x_1)$$

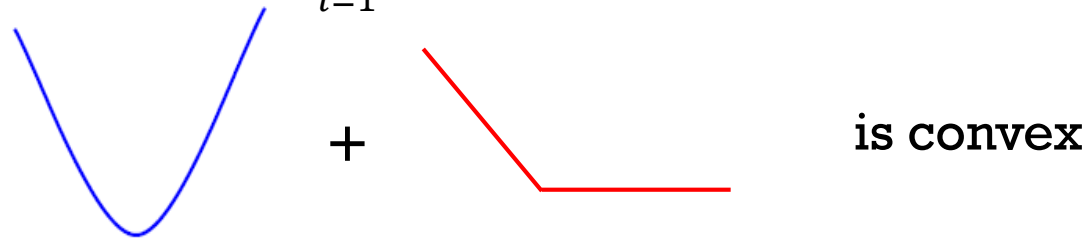
$$g((1-t)x_0 + tx_1) \leq (1-t)g(x_0) + tg(x_1)$$

$$\text{Hence } (\alpha f + \beta g)((1-t)x_0 + tx_1) \leq (1-t)(\alpha f + \beta g)(x_0) + t(\alpha f + \beta g)(x_1)$$

non-negative sum of convex functions is convex

CONVEX FUNCTION PROPERTIES

SVM's loss function: $\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b))$



→ local optimal is global optimal

→ Gradient descent always leads to global optimal regardless of initialization.

Lemma: Let Ω be a convex set in \mathbb{R}^d , f be a convex function on Ω .

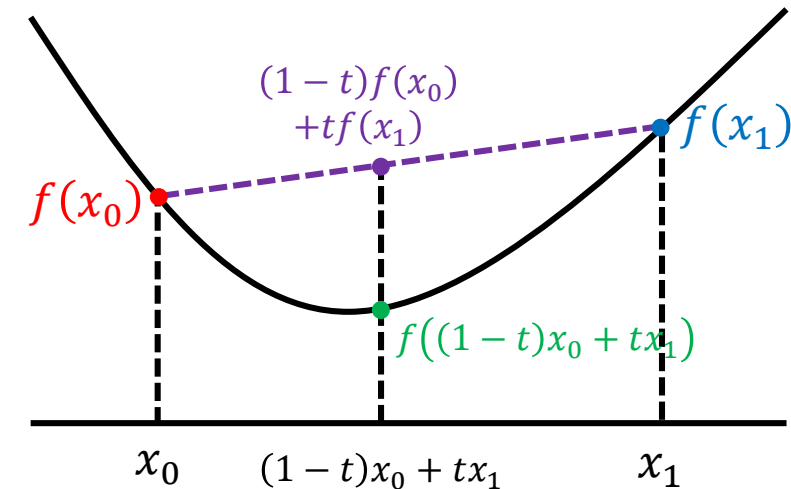
If $x_0, x_1 \in \Omega$ are local minimal points of f , then $f(x_0) = f(x_1)$.

Proof: WLOG assume $f(x_0) < f(x_1)$. For all $0 \leq t < 1$, one has

$$f((1-t)x_0 + tx_1) \leq (1-t)f(x_0) + tf(x_1) < f(x_1)$$

That is, $f(x) < f(x_1)$ for all $x \in \{(1-t)x_0 + tx_1 : 0 \leq t < 1\}$.

Hence x_1 cannot be local minimal, leading to contraction.



If the loss function is convex, then a locally

SVM 2019 Fall optimal point is globally optimal

GRADIENT DESCENT ALGORITHM FOR SVM

- To minimize a loss function $L(\mathbf{w}_t)$ we use the iterative update

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t \nabla_{\mathbf{w}} L(\mathbf{w}_t)$$

where η_t is the learning rate (at time t).

- First, rewrite the optimization problem as an average

$$L(\mathbf{w}) = \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{N} \sum_{i=1}^N \max(0, 1 - y_i f(\mathbf{x}_i))$$

where $\lambda = 1/(NC)$ and $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + \mathbf{b}$

SUB-GRADIENT DESCENT ALGORITHM FOR SVM

$$L(\mathbf{w}) = \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^N \mathcal{L}(\mathbf{x}_i, \mathbf{y}_i, \mathbf{w})$$

- The iterative update is

$$\begin{aligned} \mathbf{w}_{t+1} &\leftarrow \mathbf{w}_t - \eta_t \nabla_{\mathbf{w}} L(\mathbf{w}_t) \\ &\leftarrow \mathbf{w}_t - \eta_t \left(\lambda \mathbf{w}_t + \frac{1}{n} \sum_{i=1}^n \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{x}_i, \mathbf{y}_i, \mathbf{w}_t) \right) \end{aligned}$$

where η_t is the learning rate.

- In the Pegasos algorithm the learning rate is set at $\eta_t = \frac{1}{\lambda t}$
- Alternative: Stochastic gradient decent.

