

Gradient Descent for Support Vector Machines and Subgradients

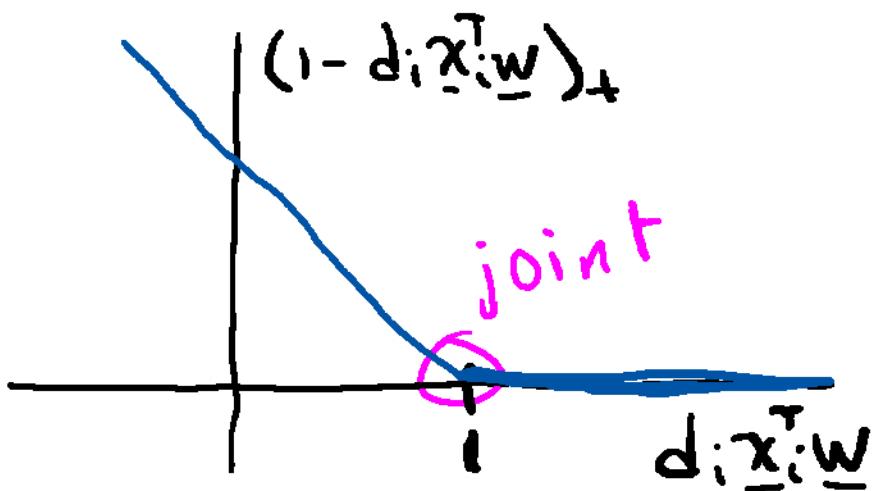
Objectives

- develop a gradient descent algorithm for SVMs
- introduce subgradients for convex but non differentiable cost functions

Support vector machines require iterative algorithms ²

$$\min_{\underline{w}} \sum_{i=1}^N (1 - d_i \underline{x}_i^\top \underline{w})_+ + \lambda \|\underline{w}\|_2^2$$

↑
labels features hinge loss regularization



No closed form solution
Convex

⇒ gradient descent

Problem: hinge loss not differentiable

★ (Not differentiable, but convex → Use a gradient-descent based solving to solve)

Subderivatives generalize derivatives

- Convex, but non-differentiable $f(x)$

★ (次导数是在处理不可微函数的凸函数时, 对传统导数的推广)
核心思想: 当 f 在某点可微 \rightarrow 次导数就是导数
不可微 \rightarrow 导数不存在, 定义次导数

Derivatives -

$$d(x_0) = \lim_{x \rightarrow x_0} \frac{f(x) - f(x_0)}{x - x_0} \quad (\text{Tangent line at } x_0)$$

Convex:

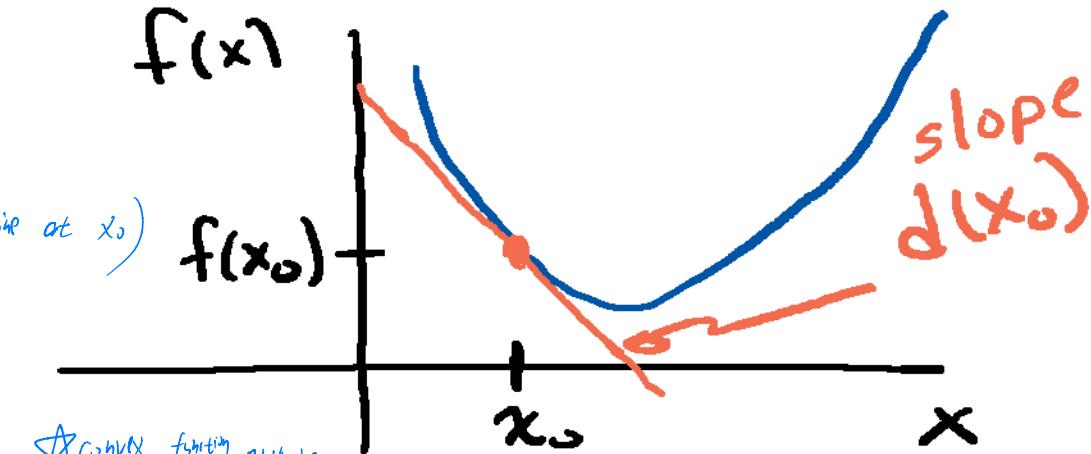
$$f(x) \geq f(x_0) + d(x_0)(x - x_0)$$

Subderivative (convex)

$$\text{Any } d_s(x_0): f(x) \geq f(x_0) + d_s(x_0)(x - x_0)$$

$$x < 1: d_s(x) = -\frac{1}{2}; \quad x > 1: d_s(x) = \frac{1}{2}$$

★ 次导数是个数, 而是满足 $f(x) \geq f(x_0) + d_s(x_0)(x - x_0)$ 的所有斜率的集合

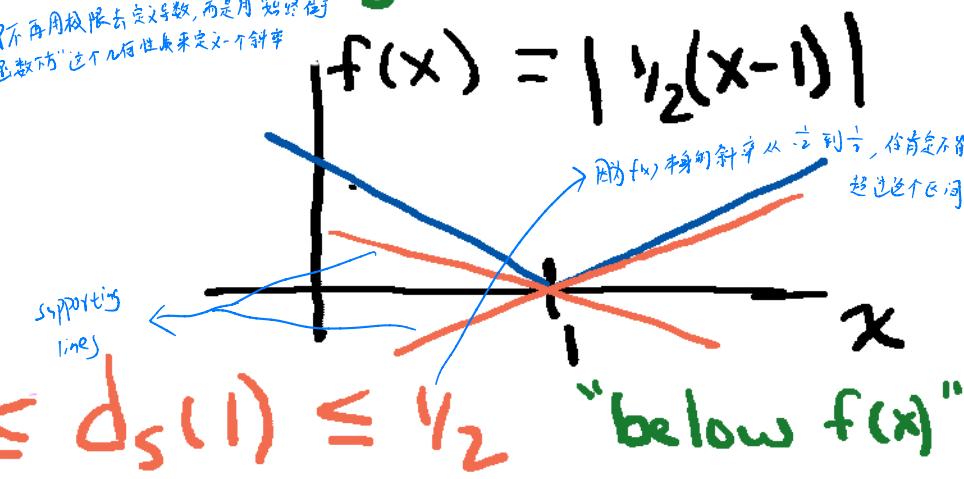


★ convex function must be

"above tangent line"

★ 不再用极限去定义导数, 而是用“切线斜率”这个几何性质来定义一个斜率

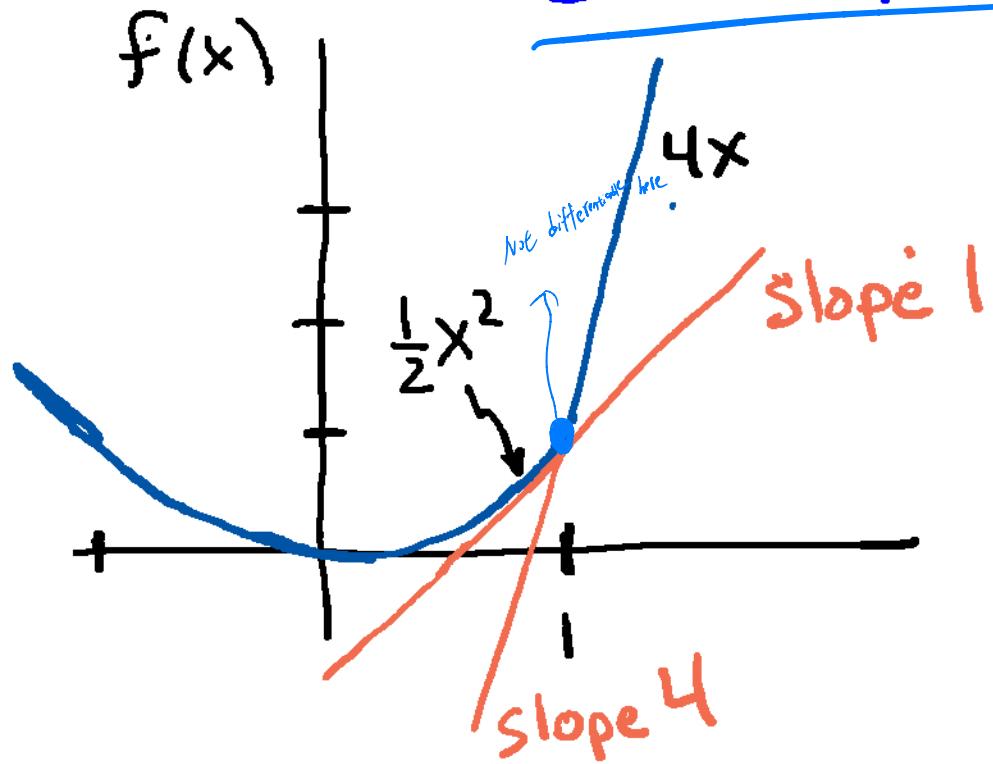
$$f(x) = |\frac{1}{2}(x-1)|$$



$$-\frac{1}{2} \leq d_s(1) \leq \frac{1}{2}$$

"below $f(x)$ "

Sub derivatives produce "reasonable" downhill directions



Example: $f(x) = \begin{cases} \frac{1}{2}x^2 & x < 1 \\ 4x & x \geq 1 \end{cases}$

convex

Subderivative

$$d_s(x) = \begin{cases} x > x < 1 \\ 4 & x > 1 \\ [1, 4] & x = 1 \end{cases}$$

$\stackrel{=\frac{d}{dx}(\frac{1}{2}x^2)=x}{\stackrel{=\frac{d}{dx}(4x)=4}{\downarrow}}$

(Any interval between 1 and 4)

Subgradients generalize gradients

- Convex,

nondifferentiable $l(\underline{w})$

Gradients -

$$l(\underline{w}) \geq l(\underline{w}_0) + (\underline{w} - \underline{w}_0)^T \underline{V}(\underline{w}_0) \quad \underline{V}(\underline{w}) = \nabla_{\underline{w}} l(\underline{w})$$

★ "above tangent plane" ($\sum_{i=1}^m (w_i - w_{0i}) \frac{d}{dw_i} l(\underline{w})$)

Subgradients -

$$\text{Any } \underline{V}(\underline{w}) : l(\underline{w}) \geq l(\underline{w}_0) + (\underline{w} - \underline{w}_0)^T \underline{V}(\underline{w}_0)$$

★ 在 w_0 点，任一个向量 $V(w_0)$ ，如果由它定义的“支撑平面” $y = l(w_0) + (\underline{w} - \underline{w}_0)^T V(w_0)$ 始终低于 $l(w)$ ，那么这个 $V(w_0)$ 就是 w_0 的一个 subgradient

★ Gradient descent optimization: replace gradient with subgradient

就是 subderivative 在多维中的叫法
(即有 "supporting plane" 的梯度向量的集合)

在 w_0 点的梯度向量 $\left[\frac{\partial L}{\partial w_1} \frac{\partial L}{\partial w_2} \dots \right]$

Tangent plane

Gradient descent for SVMs

$$\ell(\underline{w}) = \sum_{i=1}^N (1 - d_i \underline{x}_i^\top \underline{w})_+ \quad \rightarrow \text{subgradient}$$

↑ Hinge Loss

$$\ell_i(\underline{w}) = (1 - d_i \underline{x}_i^\top \underline{w})_+ = \begin{cases} 1 - d_i \underline{x}_i^\top \underline{w} & d_i \underline{x}_i^\top \underline{w} < 1 \\ 0 & d_i \underline{x}_i^\top \underline{w} \geq 1 \end{cases}$$

↓ For each sample i

Subgradient

$$\nabla_i(\underline{w}) = \begin{cases} -d_i \underline{x}_i & d_i \underline{x}_i^\top \underline{w} < 1 \\ 0 & d_i \underline{x}_i^\top \underline{w} \geq 1 \end{cases} = -d_i \underline{x}_i I_{\{d_i \underline{x}_i^\top \underline{w} < 1\}}$$

↑ indicator function

$$\text{Cost } f(\underline{w}) = \ell(\underline{w}) + \lambda \|\underline{w}\|_2^2$$

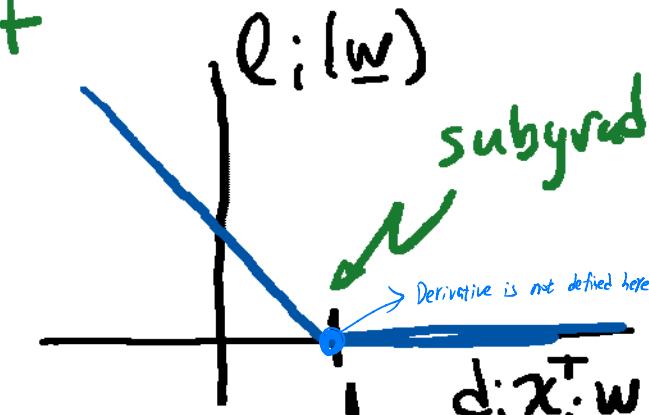
$$\Rightarrow \nabla f(\underline{w})|_{\underline{w}^{(k)}} = \sum_{i=1}^N (-d_i \underline{x}_i I_{\{d_i \underline{x}_i^\top \underline{w}^{(k)} < 1\}}) + 2\lambda \underline{w}^{(k)}$$

↑ Gradient of $\lambda \|\underline{w}\|_2^2$

Gradient descent

$$\underline{w}^{(k+1)} = \underline{w}^{(k)} - \tau \nabla f(\underline{w})|_{\underline{w}^{(k)}}$$

↑ 最陡的负向, 下坡



**Copyright 2019
Barry Van Veen**