

Gradient Descent for Support Vector Machines and Subgradients

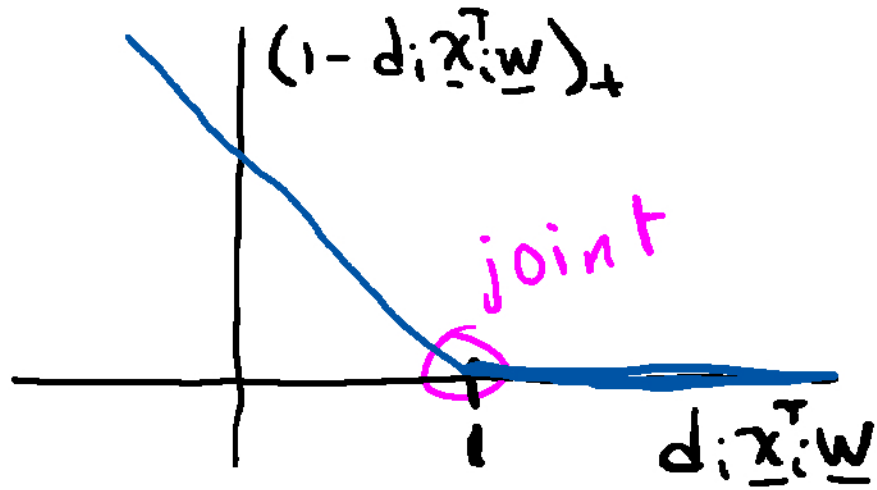
Objectives

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

Support vector machines require iterative algorithms ²

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

labels features hinge loss regularization



No closed form solution

Convex

\Rightarrow gradient descent

Problem: hinge loss not differentiable

Subderivatives generalize derivatives 3

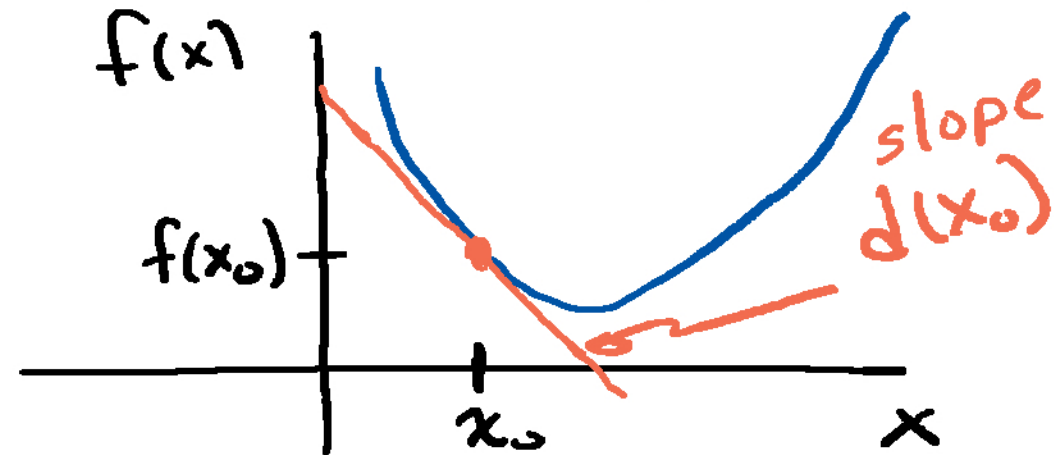
– Convex, but nondifferentiable $f(x)$

Derivatives –

$$d(x_0) = \lim_{x \rightarrow x_0} \frac{f(x) - f(x_0)}{x - x_0}$$

Convex:

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

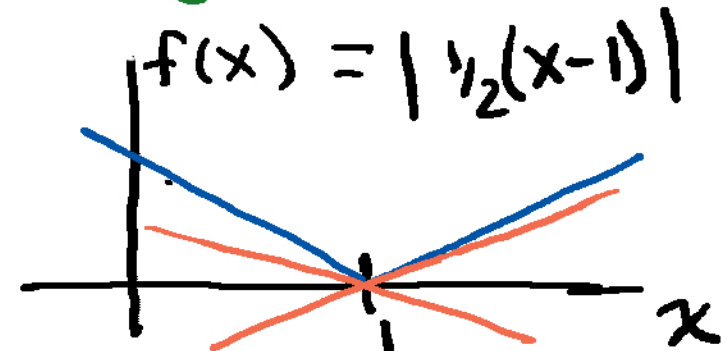


"above tangent line"

Subderivative (convex)

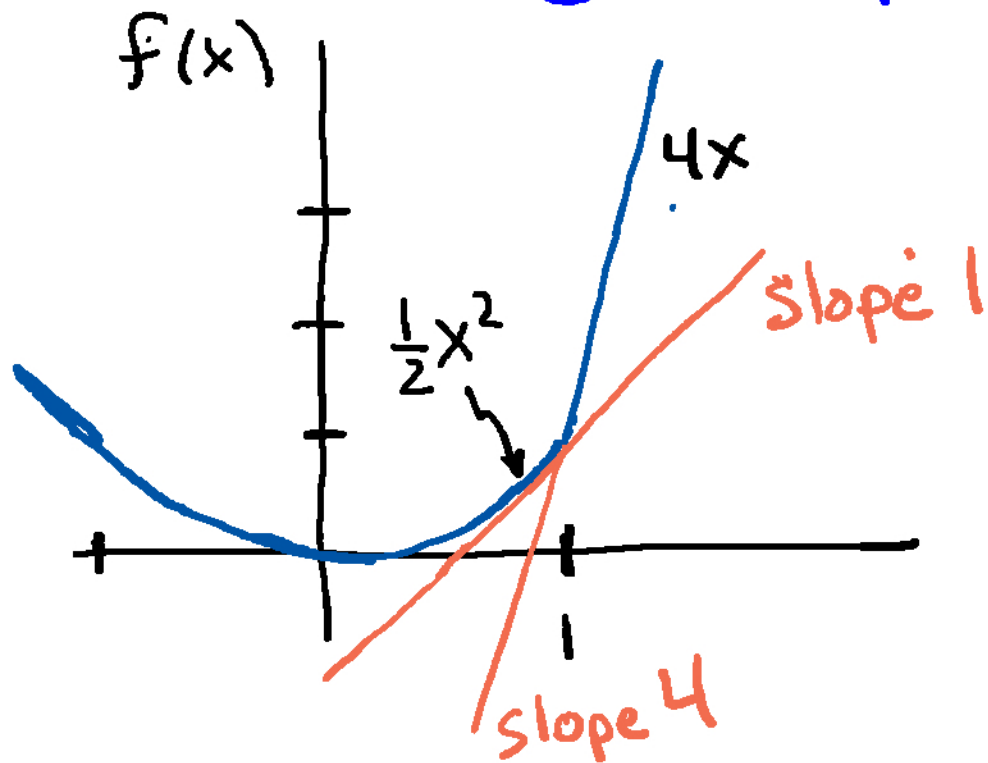
Any $d_s(x_0): f(x) \geq f(x_0) + d_s(x_0)(x - x_0)$

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



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

Sub derivatives produce "reasonable"
downhill directions 4



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

Subderivative

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

Subgradients generalize gradients 5
- 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" $\left(\sum_{i=1}^M (w_i - w_{0i}) \frac{d}{dw_i} l(\underline{w}) \right)$

Subgradients -

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

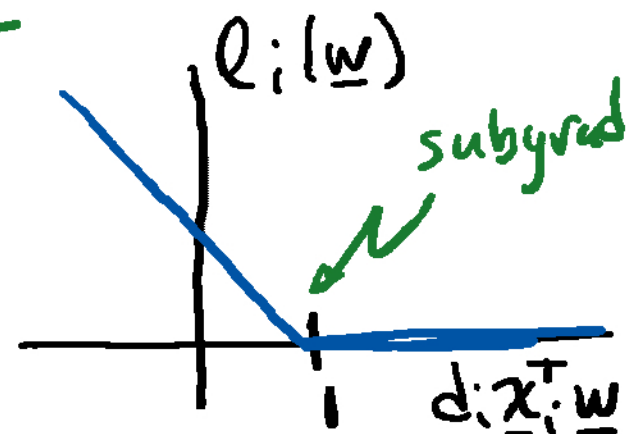
Gradient descent optimization: replace
gradient with subgradient

Gradient descent for SVMs

6

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

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



Subgradient

$$\underline{v}_i(\underline{w}) = \begin{cases} -d_i \underline{x}_i & d_i \underline{x}_i^T \underline{w} < 1 \\ 0 & d_i \underline{x}_i^T \underline{w} \geq 1 \end{cases} = -d_i \underline{x}_i \mathbb{I}_{\{d_i \underline{x}_i^T \underline{w} < 1\}}$$

indicator function

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

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

Gradient descent

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

**Copyright 2019
Barry Van Veen**