

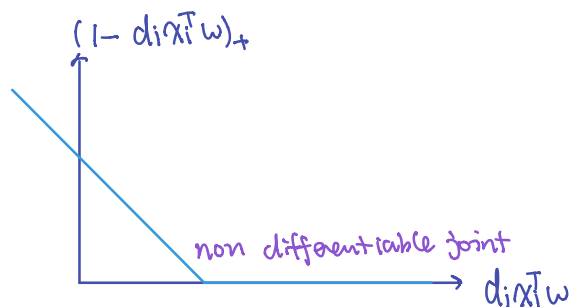
Support vector machines require iterative algorithms

$$\min_w \sum_{i=1}^N (1 - d_i x_i^T w)_+ + \lambda \|w\|_2^2$$

\uparrow labels \uparrow features \uparrow hinge loss \uparrow regularization

No closed form solution.

but convex function \Rightarrow gradient descent



Problem: hinge loss non-differentiable

Subderivatives generalize derivatives

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

Derivatives -

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

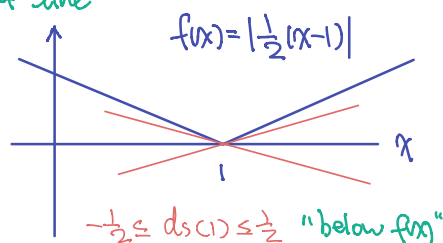
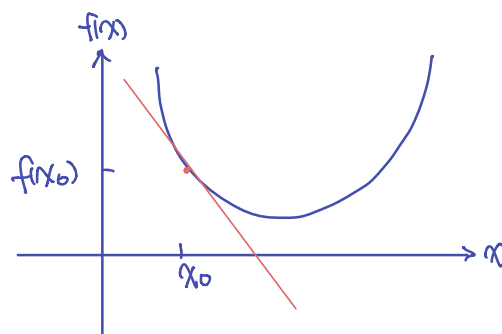
For a convex function:

$$f(x) \geq f(x_0) + d(x_0)(x - x_0) \quad \text{"above tangent line"}$$

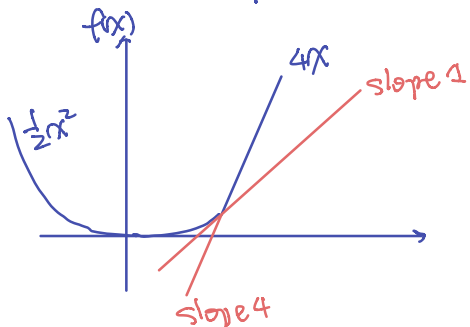
Subderivative (convex)

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

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



Subderivatives produce "reasonable" downhill directions



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

- Convex, non differentiable $l(x)$

Gradients-

$$l(w) \geq l(w_0) + (w - w_0)^T \nabla l(w_0) \quad , \quad \nabla l(w) = \nabla_w l(w)$$

$$\text{"above tangent plane"} \quad \left(\sum_{i=1}^M (w_i - w_{0i}) \frac{d}{dw_i} l(w_0) \right)$$

Subgradients-

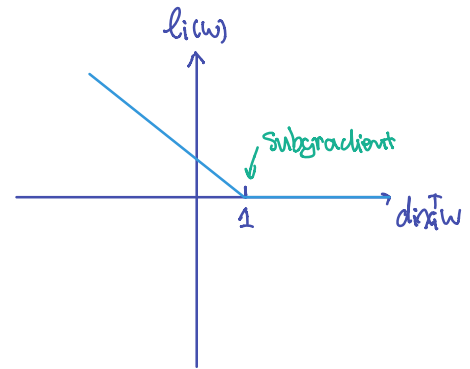
$$\text{Any } v(w): l(w) \geq l(w_0) + (w - w_0)^T v(w_0)$$

Gradient descent optimization: replace gradient with subgradient

Gradient descent for SVMs

$$l(w) = \sum_{i=1}^N (1 - d_i x_i^T w)_+ \rightarrow \text{subgradient}$$

$$l_i(w) = (1 - d_i x_i^T w)_+ = \begin{cases} 1 - d_i x_i^T w & d_i x_i^T w < 1 \\ 0 & d_i x_i^T w \geq 1 \end{cases}$$



Subgradient

$$v_i(w) = \begin{cases} -d_i x_i & d_i x_i^T w < 1 \\ 0 & d_i x_i^T w \geq 1 \end{cases} = -d_i x_i I_{\{d_i x_i^T w < 1\}}$$

Indicator function

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

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

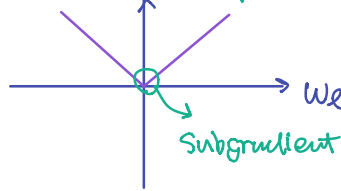
Gradient descent

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

Example: Gradient Descent for LASSO

$$f(w) = \sum_{i=1}^N (d_i - x_i^T w)^2 + \lambda \|w\|_1 = \sum_{i=1}^N \left\{ \underbrace{(d_i - x_i^T w)^2}_{|w_e|} + \underbrace{\frac{\lambda}{N} \|w\|_1}_{f_i(w)} \right\}$$

Consider $\nabla w = \sum_{e=1}^M |w_e|$



$$\frac{d}{dw_e} |w_e| = \begin{cases} \text{sign}(w_e) & w_e \neq 0 \\ [-1, 1] & w_e = 0 \end{cases}$$

↑
"0" popular

write $\nabla \|w\|_1 = \text{sign}(w)$

$$\nabla_w f_i(w) = -2(d_i - x_i^T w) x_i + \frac{\lambda}{N} \text{sign}(w)$$

$$w^{(k+1)} = \tau (d_{ik} - x_{ik}^T w^{(k)}) x_{ik} - \frac{\lambda \tau}{2N} \text{sign}(w^{(k)})$$