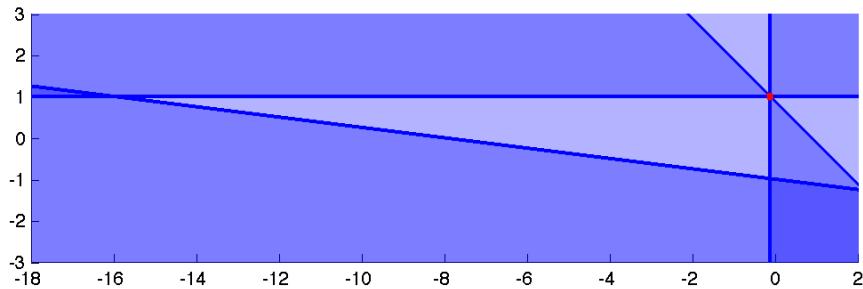
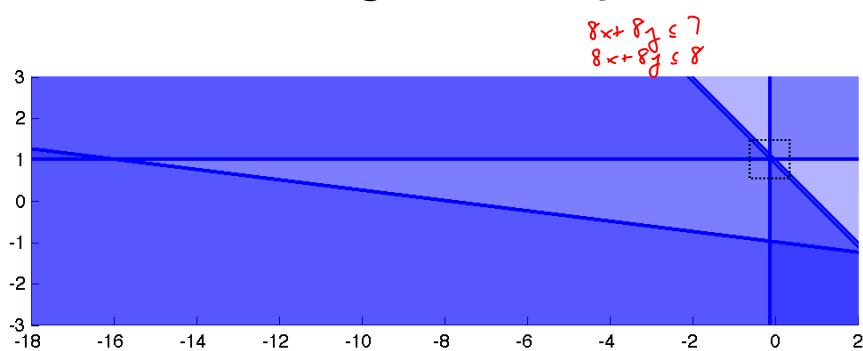


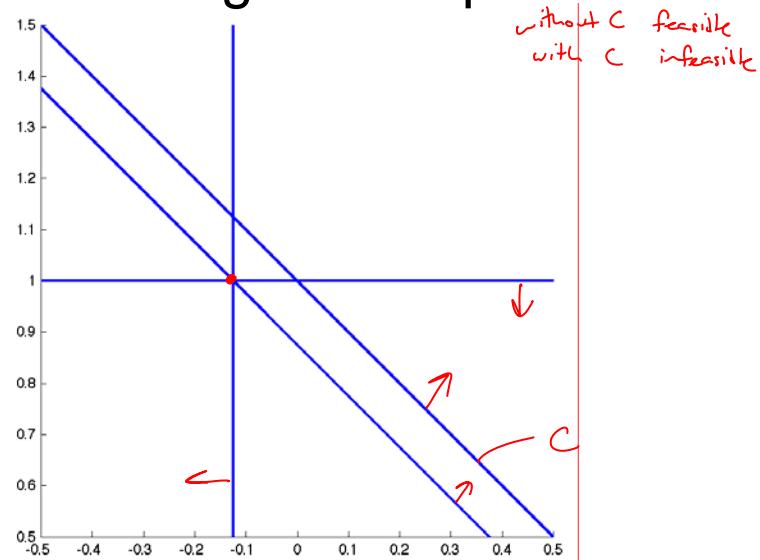
## Bit length example



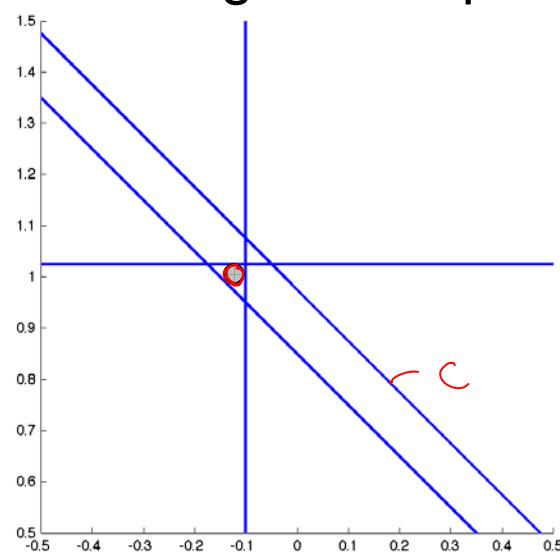
## Bit length example



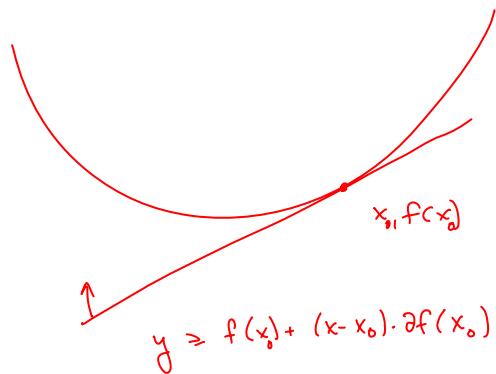
## Bit length example



## Bit length example



## What's a subgradient?



## Subgradients for SVMs

- $\min_w L(w) = \|\underline{w}\|^2 + (C/m) \sum_i h(-y_i x_i^T w)$

- $h(z) = \max \{0, 1+z\}$

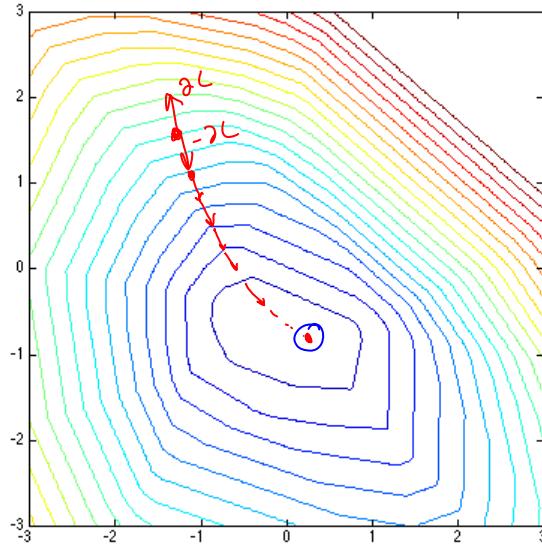
- Subgradient of  $h(z)$ :

$$\partial h(z) = \begin{cases} 0 & z < -1 \\ 1 & z > -1 \\ [0,1] & z = -1 \end{cases}$$

- Subgradient of  $L(w)$  wrt  $w$ :

$$\underline{\partial L(w)} = \underline{2w} + \frac{1}{m} \sum_i \underline{\partial h(-y_i x_i^T w)} \cdot \underline{(-y_i x_i)}$$

## Subgradient descent



## Subgradient descent

Start w/  $\underline{x}_0$

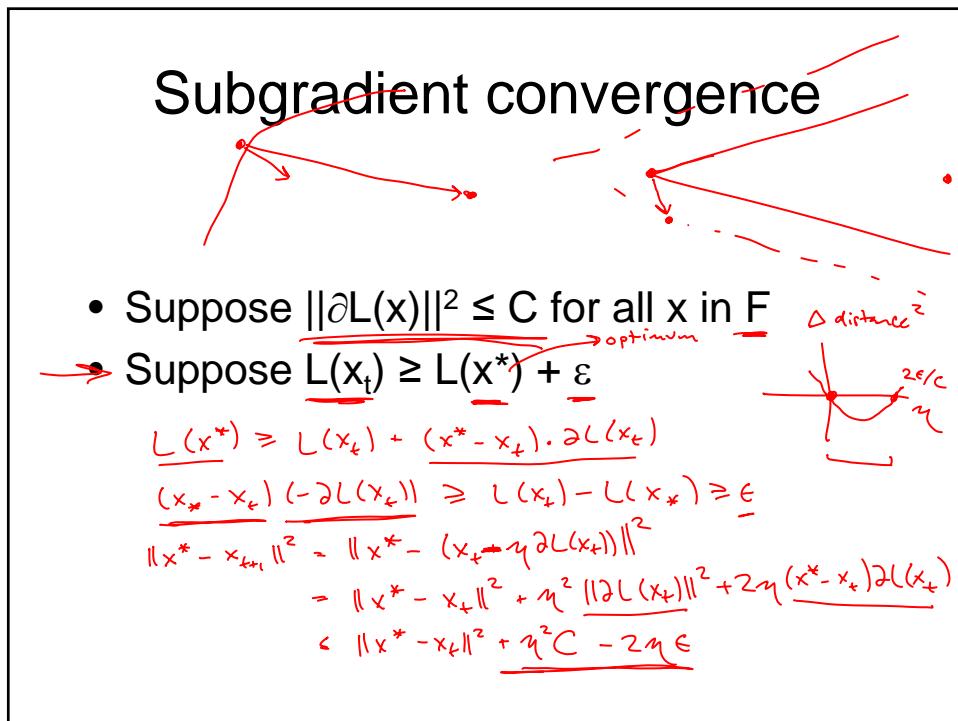
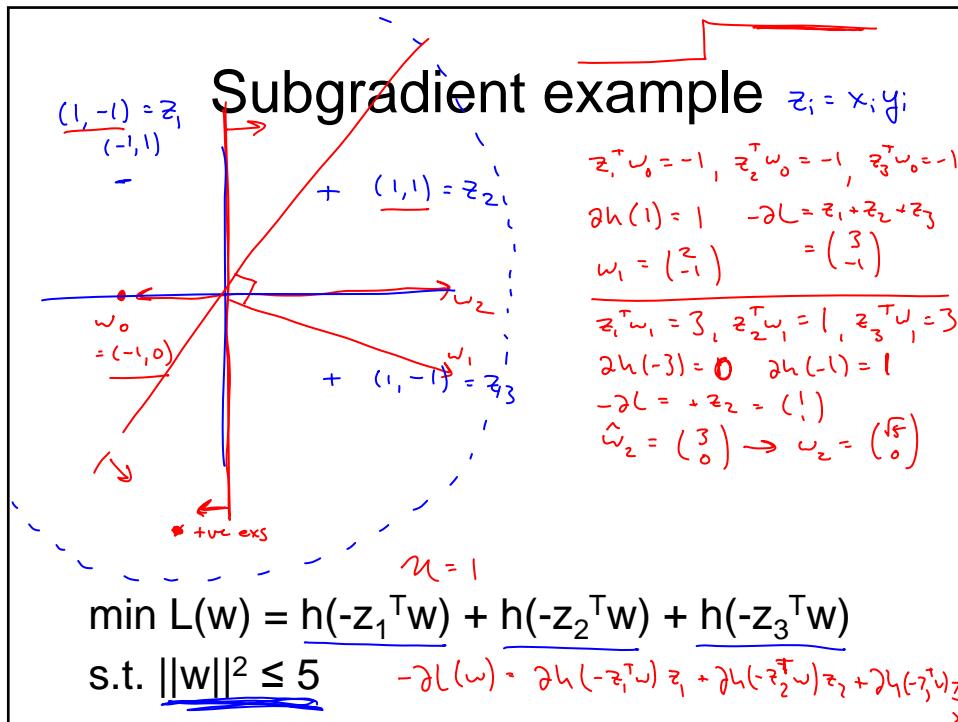
- While not tired:  $\eta_t$  = learning rate

$$\underline{g}_t = \text{(estimate of)} \quad \underline{\partial f(x_t)}$$

$$\underline{x}_{t+1} = \underline{x}_t - \eta_t \underline{g}_t$$

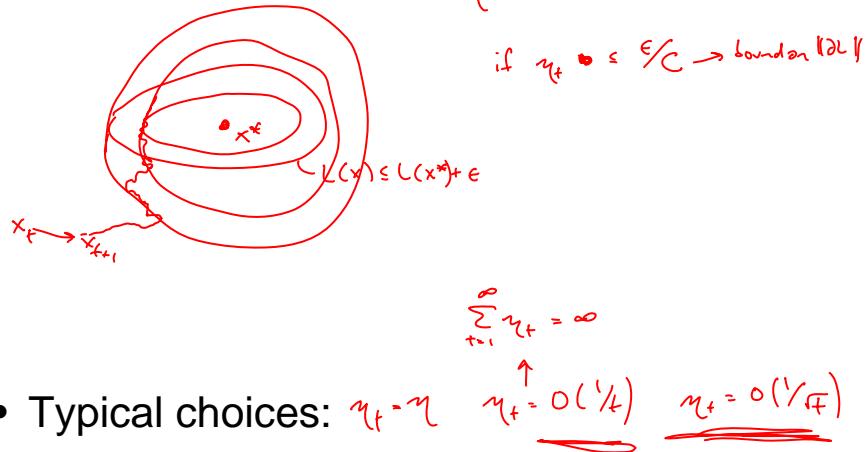
$$\underline{x}_{t+1} := \underline{\bigcap_F x_{t+1}}$$

projection onto feasible region F



## Setting step size

- If we knew  $\varepsilon$ , could set good step size  $\eta$
- But we don't! So:  $\eta_t \rightarrow 0$  as  $t \rightarrow \infty$



## Stochastic subgradient

- In SVM (and many other ML problems),  $L(w)$  contains big sum of simple terms
$$\min_w L(w) = \|w\|^2 + (\underbrace{C/m}_{\leq C\|x_i\|} \sum_i h(-y_i x_i^T w))$$

$$\partial L(w) = 2w - \underbrace{\sum_i \partial h(-y_i x_i^T w)}_{\text{norm of term } i} (y_i x_i) \leq \|x_i\|$$
- Approximate sum by sampling terms

$$\|\partial_i\| \leq \frac{C \partial h(-y_i x_i^T w)}{C\|x_i\|} y_i x_i \quad \partial_S = 2w - \frac{C}{\|S\|} \sum_{i \in S} \partial_i \quad E(\partial_S) = \partial L(w)$$

$$E(-\partial_S^T (x - x^*)) = -\partial L(w) \cdot (x - x^*)$$

$$\text{S random, } |S| = k: \text{Var}(\eta \partial_S) \leq \eta^2 \frac{C^2 \|x\|^2}{k} = O(\eta^2)$$

# When do we stop?

$f(x)$

for SVMs  
get accuracy  $\in$   
 $O(1/\epsilon)$

- Feasible region diameter  $\|F\|$

$$\begin{aligned} f(x_t) \geq \underline{f(x^*)} &\geq f(x_t) + (x^* - x_t) \cdot \nabla f(x_t) \geq f(x_t) - \underline{\|x^* - x_t\| \|\nabla f(x)\|} \\ &\geq f(x_t) - \underline{\|F\| \|\nabla f(x_t)\|} \end{aligned}$$

- Typical ML generalization bound:

$$\begin{aligned} \mathbb{E}_w L(\text{new ex}, w) &\leq L(\text{train}, w) + \underline{\text{stuff}} \\ w^* = \text{opt} \quad L(\text{train}, w') &\leq L(\text{train}, w^*) + O(\text{stuff}) \\ &\in (L(\text{new ex}, w')) \leq L(\text{train}, w') + \text{stuff} \\ &\leq L(\text{train}, w^*) + \underline{O(\text{stuff})} + \text{stuff} \xrightarrow{O(\text{stuff})} \end{aligned}$$