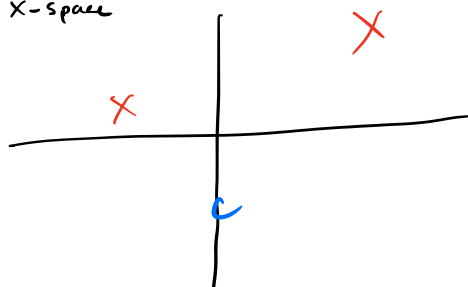


Gradient Descent, Stochastic gradient descent, Perceptron learning algorithm. Feature vs. weight space. Max margin classifier, aka hard margin SVM.

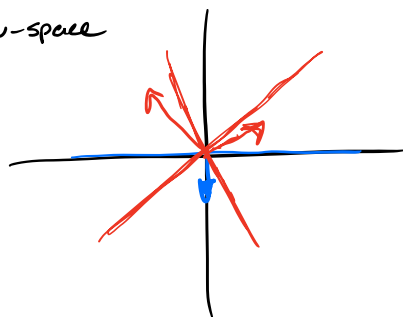
Perceptron continued

$$R(w) = \sum_{i \in V} -y_i x_i^T w \quad \forall \{x_i : y_i x_i^T w < 0\}$$

X-space



w-space



Optimization algorithm: Gradient descent on R :

- Starting point w (not origin), find gradient w.r.t. w
 \hookrightarrow step in opposite direction to steepest ascent

$$R(w) = \sum_{i \in V} -y_i x_i^T w = - \sum_{i \in V} y_i x_i$$

$w \leftarrow$ non zero starting point (use a sample point)

while $R(w) > 0$:

$$V \leftarrow \{i \text{'s for } y_i x_i^T w < 0\}$$

$$w \leftarrow w + \epsilon \sum_{i \in V} y_i x_i$$

$\epsilon > 0$ is the step size aka learning rate

slow! $O(nd)$ time !!!

Optimization alg. 2: Stochastic gradient descent

- each step, pick one misclassified X_i
- do gradient descent on loss function of $L(X_i^T w, y_i)$
- "perceptron alg" $O(d)$ time

while true $y_i X_i^T w < 0$

$$w \leftarrow w + \epsilon y_i X_i$$

return w

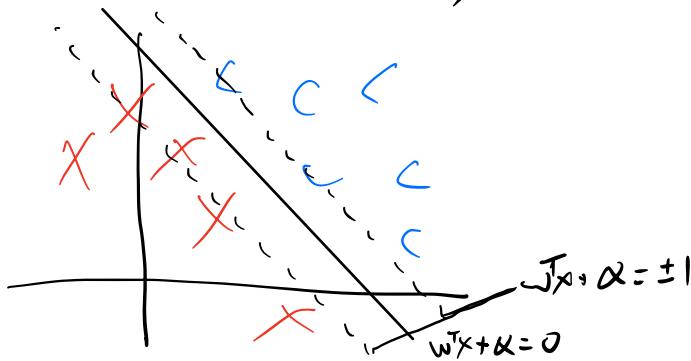
what if $\alpha \neq 0$

Add a fictitious dimension: \mathbb{R}^{d+1} $[X_i]_{d+1} = 1$

$$f(x) = w^T x + \alpha = [w \ \alpha] \begin{bmatrix} x \\ 1 \end{bmatrix}$$

MAX MARGIN CLASSIFIER

margin = dist from decision boundary to closest point



$$y_i (w^T x_i + \alpha) \geq 1 \text{ for } i \in [1, n]$$

$\|w\|_2 = 1$, signed dist to X_i is $w^T x_i + \alpha$

otherwise, $\frac{w^T}{\|w\|_2} x_i + \frac{\alpha}{\|w\|_2}$

margin is $\min_i \frac{1}{\|w\|_2} |w^T x_i + \alpha| \geq \frac{1}{\|w\|_2}$

To max margin, min w :

QP $\rightarrow \min_w \|w\|_2^2 : y_i (x_i^T w + \alpha) \geq 1 \text{ for } i=1 \dots n$

"hard margin SVM" margin $\frac{1}{\|w\|_2}$ slack $\frac{2}{\|w\|_2}$