Perspective & connection to DS/ML:

In data science / machine learning we often need to optimize functions of the form

$$\min_{w} \frac{1}{N} \sum_{i=1}^{N} f\left(w ; (x_i, y_i)\right) + \lambda R(w)$$

model parameters — loss function — training examples or data — regularization

Examples / applications:

× linear regression
× logistic regression
# least squares fit to a model
* PCA
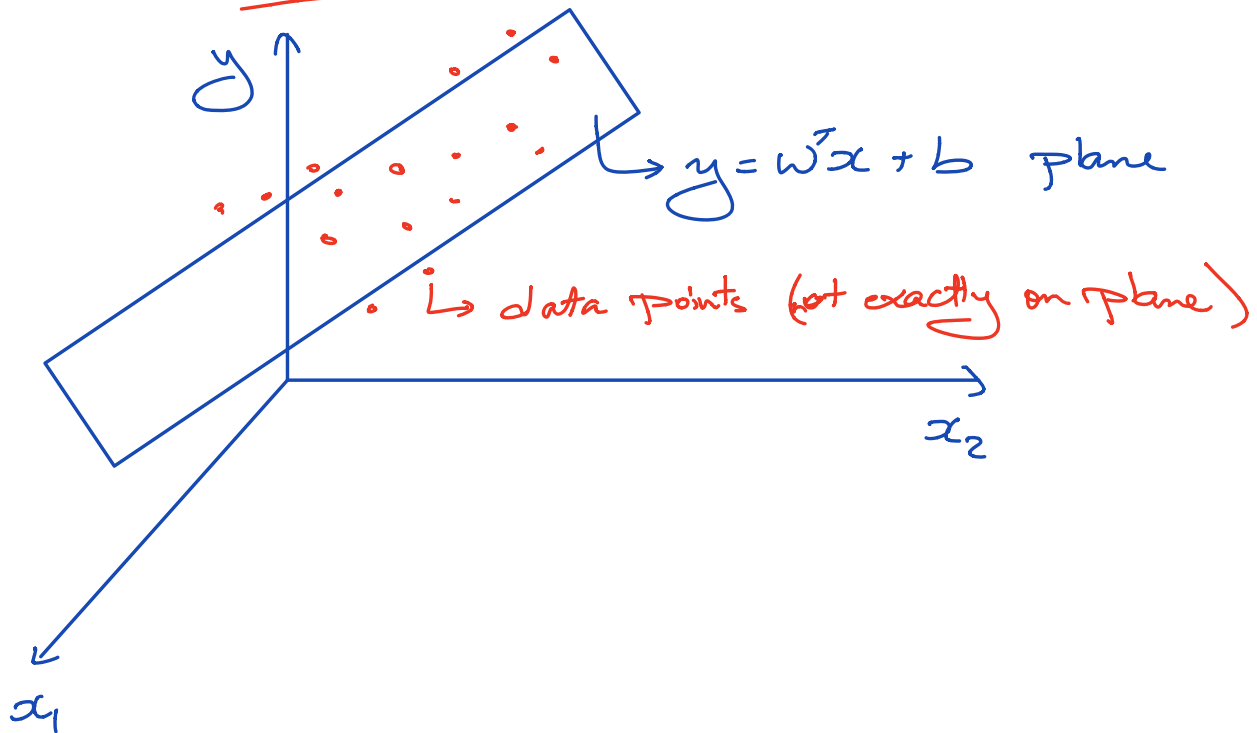* Neural network loss
× Support vector machines
*(K-means) clustering

° ° °

× **Linear regression:**

Given data $(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$, $i = 1, \ldots, N$

"regressors" ↗ ↪ dependent variable

linear regression assumes that the relationship between $y_i$ & $x_i$ is linear, up to an error or noise term so

$J \in \mathbb{R}$ unknown

$$y_i = w^T x_i + b + \varepsilon_i$$

↑ ↑ ↑
$\in \mathbb{R}^d$  $\in \mathbb{R}^d$  $\in \mathbb{R}$, noise, <u>unknown</u>
<u>unknown</u>



$y = w^T x + b$ plane

↪ data points (not exactly on plane)

<u>Goal</u> : estimate $w$ & $b$

e.g. $\varepsilon_i$ is Gaussian random variable

under some assumptions, it makes sense to solve

$$\min_{w,b} \frac{1}{N} \sum_{i=1}^{N} |w^T x_i + b - y_i|^2$$

( linear least squares, has closed form solution )

Sometimes we would like our model to be "parsimonious", i.e., we would like most entries of $w$ to be zero, so we introduce the regularizer

$$R(w) = \|w\|_1 = \sum_{i=1}^{d} |w_i|$$

and we solve

$$\min_{w,b} \frac{1}{N} \sum_{i=1}^{N} |w^T x_i + b - y_i|^2 + \lambda \|w\|_1$$

$\hookrightarrow$ "Lasso", no closed form sol'n

not differentiable

\* **Logistic regression** (Classification)

Data : $(x_i, y_i) \in \mathbb{R}^d \times \{-1, 1\}$ , $i = 1, \cdots, N$

logistic regression assumes that the
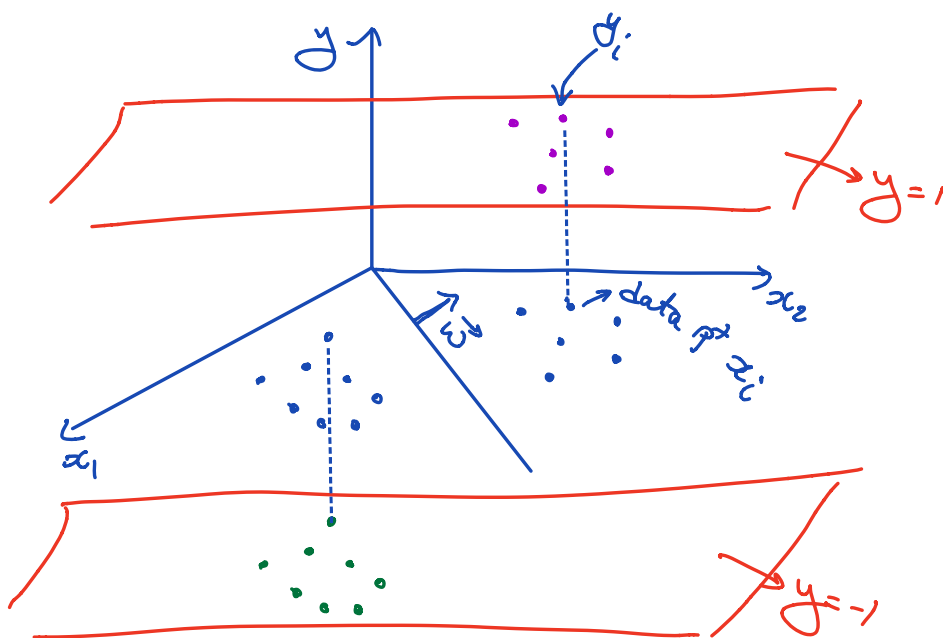relationship between $y_i$ & $x_i$ is

$$y_i = \text{sign}(w^T x_i + b + \mathcal{E}_i)$$
$$= \text{sign}(\tilde{w}^T \tilde{x}_i + \mathcal{E}_i)$$

"noise", "error"

$\quad\hookrightarrow$ CDF

$$\frac{1}{1 + e^{-z}}$$

where $\quad \tilde{w} = (w, b) \in \mathbb{R}^{d+1}$
$$\tilde{x} = (x, 1) \in \mathbb{R}^{d+1}$$
So $\quad w^T x + b = \tilde{w}^T \tilde{x}$

Subject of a lot of our
assignments

$$\min_{w} \quad \frac{1}{N} \sum_{i=1}^{N} \log \left( 1 + e^{-w^T x_i \cdot y_i} \right)$$

\* Least squares fit to a model

Example: HW 2 where we
fit the parameters of
a hyperbola based on
data.

\* Nonlinear least squares:

Data: $(x_i, y_i)$ as before

Model: $y_i = h_w(x_i)$
$\qquad\qquad \llcorner$ nonlinear function
$\qquad\qquad\qquad$ parametrized by $w$

Example: $y = b e^{w^T x}$

$$\min_{\substack{b \in \mathbb{R} \\ w \in \mathbb{R}^d}} \quad \frac{1}{N} \sum_{i=1}^{N} |y_i - b \, e^{w^T x_i}|^2$$

Can lead to non-convex opt. problems

(training) <u>Neural Networks</u>:

Neural nets are functions of the form

$$y = g_L \circ g_{L-1} \circ \cdots \circ g_1 (x) =: G_{w,b}(x)$$

↑ composition

$$= g_L(g_{L-1}( \cdots g_2(g_1(x))))$$

where $g_j(z) = \rho(W^{(j)} z + b^{(j)})$

↗ non-linear ↑ matrix ↘ vector
Guntion

e.g. relu $\rho(z) = \max(0, z)$

Given data $(x_i, y_i)$, $i = 1, \cdots, N$

want to solve for $(W^{(j)}, b^{(j)})$, $j = 1, \cdots, L$

$\Rightarrow$ Solve

$$\min_{\substack{w^{(j)},\, b^{(j)} \\ j=1,\dots,L}} \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}\left( (w^{(j)}, b^{(j)}); (x_i, y_i) \right)$$
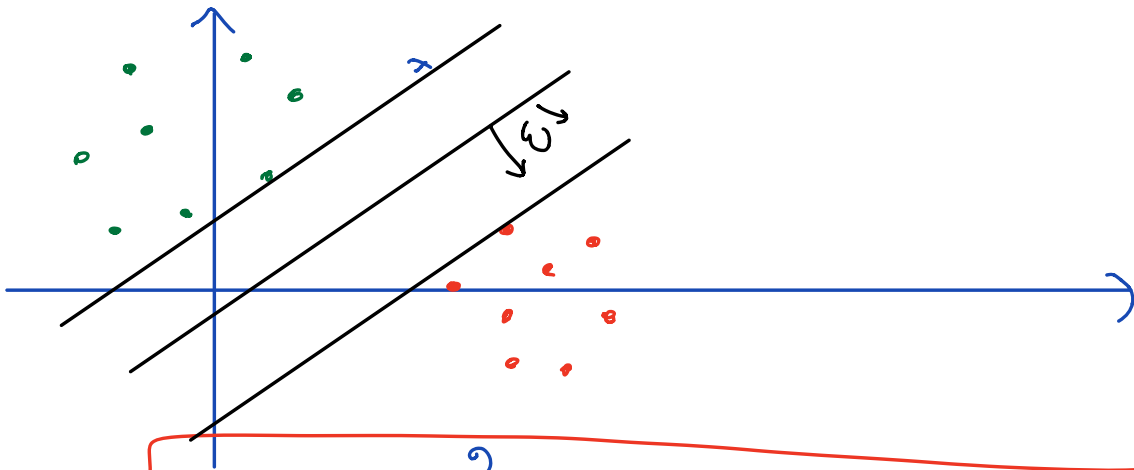
$\underset{\uparrow}{\phantom{x}}$ some loss function

No closed form, can be hard to optimize even with methods from this course

$\Rightarrow$ Stochastic GD
(17.3 B)

## Support Vector Machines (Classification)

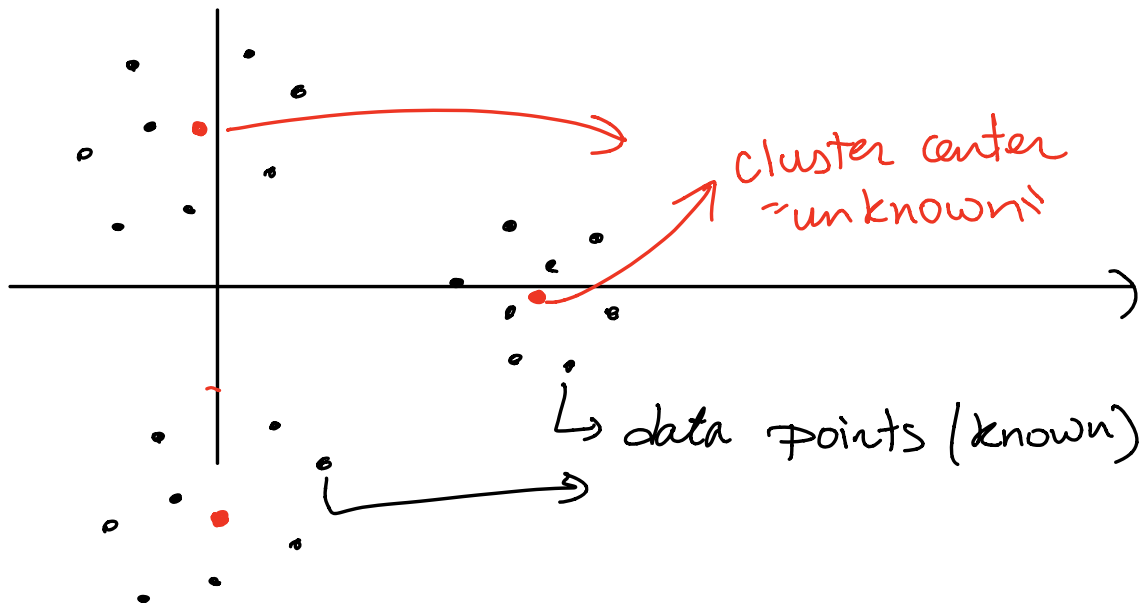__Data__ : $(x_i, y_i) \in \mathbb{R}^d \times \{-1, 1\}$, $i = 1, \dots, N$



$$\max_{w,b} \frac{2}{\|w\|}$$

$$\text{s.t.} \begin{cases} \langle w, x_i \rangle + b \geq 1 & \forall \text{ red point} \\ \langle w, x_j \rangle + b \leq -1 & \forall \text{ green point} \end{cases}$$

## K-means Clustering :

Given data $x_i \in \mathbb{R}^d$ (no $y_i$'s)

cluster center
"unknown"

$\hookrightarrow$ data points (known)

the goal is to "cluster" the data
into groups where $x_i$'s in the
same cluster are closer to their
cluster "center" than to other
cluster centers.

$k$ clusters

$$\text{minimize} \atop S \qquad \sum_{i=1}^{k} \sum_{x \in S_i} \|x - \mu_i\|^2$$

$k \to$ number of clusters

$\mu_i \to$ cluster centers

$S_i \to$ the clusters (sets of $x$'s)

$\to$ turns out to be NP-hard (tough to find algorithms guaranteed to solve this)

Lots of good heuristics algorithms.

$\circ \; \circ \; \circ$