# Lecture 22: L-homogeneous KKT points & project info

Last time: minimum norm predictors
(care because: generalize well).

local plan:

* GF on exponential or logistic loss
  $\xrightarrow{t\to\infty}$ minimum norm (classification) sdn.

$$\min \frac{1}{2}\|w\|_2^2$$
$$\text{s.t. } y_i x_i^T w \geq 1 \quad \forall i$$

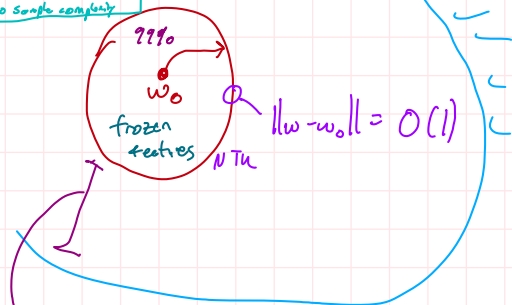$$\max \min_i y_i x_i^T w$$
$$\text{s.t. } \|w\| \leq 1$$

* closest theorem L-homogeneous preds

global plan:

| Mean-field: ① tiny int |
| ② small step ③ log growth |
| ④ no sample complexity |

↑ linear pred

??? $w_0$
$w_0$
frozen features

$\|w - w_0\| = O(1)$

NTW

asymptotic regime

$\|w_t\| \gg \|w_0\|$

good news:
* includes feature learning

bad news
* slow convergence, hard to prove / unclear what's true

intermediate regime:
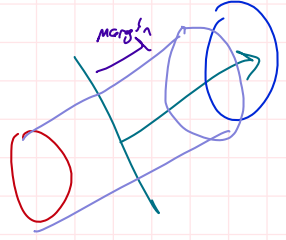* still has some randomness
* little bit of feature learning
* matches practice

# KKT points of L-homogeneous

Recall L-homogeneous predictor $F(x; cw) = c^L F(x; w)$     $c \geq 0$
(generalizes L-layer ReLU & leaky ReLU networks).

Linear case, margin

$$\max_{w \in \mathbb{R}^p} \frac{\min_i y_i F(x_i; w)}{\|w\|}$$

$$= \max_{w \in \mathbb{R}^p} \frac{\min_i \|w\|^L y_i F(x_i; \frac{w}{\|w\|})}{\|w\|}$$

$$= \max_{w \in \mathbb{R}^p} \|w\|^{L-1} \underbrace{\min_i y_i F(x_i; w/\|w\|)}_{\text{in linear case, } = 1; \text{ in } L \geq 2 \text{ case, bogus defn.}}$$

Obvious fix:

$$\max_{w \in \mathbb{R}^p} \frac{\min_i y_i F(x_i; w)}{\|w\|^L} = \max_{\|w\| \leq 1} \min_i y_i F(x_i; w).$$

**Remark:** Modern networks are not L-homogeneous (e.g., softmax, biases).

Consider two opt problems

$$\min \frac{1}{2}\|w\|^2$$
$$\text{s.t. } y_i F(x_i; w) \geq 1 \; \forall i. \quad \boxed{1}$$
$$w \in \mathbb{R}^p$$

$$\max_{w \in \mathbb{R}^p} \frac{-\ln \sum_i \exp(-y_i F(x_i; w))}{\|w\|^L} . \quad \boxed{2}$$

$\boxed{\text{Mean-field}}$ (teal box)

$\lceil$ **Theorem.** Suppose GF on $\mathcal{L}(w) = \sum_i \exp(-y_i F(x_i; w))$ and $\exists \tau \; \mathcal{L}(w_\tau) < 1$.

   ① $t \mapsto \frac{-\ln \mathcal{L}(w_t)}{\|w_t\|^L}$ (i.e., objective in $\boxed{2}$) is nondecreasing over $[\tau, \infty)$.
   [Lyu-Li].

   ② $\frac{w_t}{\|w_t\|} \to$ KKT point of $\boxed{1}$    ([Lyu-Li], [Ji-T.])
     (under o-minimal definability).
         $\lfloor$ rules out oscillations; e.g., can't have $\sin(\cdot)$ activations

**Proof remark:** ① note $v \mapsto \max_i v_i$ is 1-homogeneous,
     proof uses "approximate asymptotic 1-homogeneity" of $v \mapsto -\ln \sum_i \exp(-v_i)$.

     ② KKT point proof shows "alignment": $\left\langle \frac{w_t}{\|w_t\|}, \frac{-\nabla \mathcal{L}(w_t)}{\|\nabla \mathcal{L}(w_t)\|} \right\rangle \to 1.$

**Remark:** goodness: regularity outside nth.

     bad news: no better than nth.

**Remark:** People have shown convergence <u>global</u> nd margin soln,
     but it requires infinitely many assumptions. [Ch. Zof-Bach]