# THEORY OF MACHINE LEARNING

# LECTURE 22

NTK SUMMARY, REPRESENTATION LEARNING

# ANNOUNCEMENTS

- **Homework 3 due on Monday April 11**

- **Project presentations:** starting in two weeks! (~18 projects)

  - Dates: April 19, 22, 26 (5 projects /class), couple online
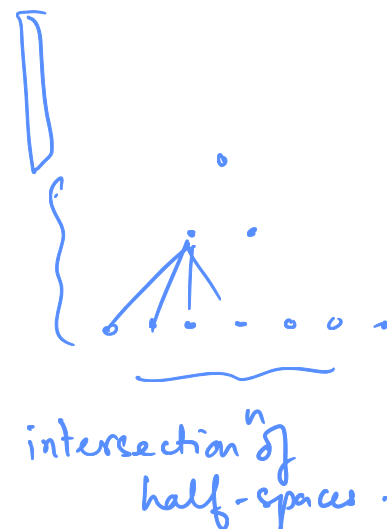
- This week and next: representation learning, robustness

**LAST WEEK**

$\langle x, \cdot \rangle$

$f(A^{(1)}) = b_1,$
$f(A^{(2)}) = b_2 \quad \cdots$

$A^{(1)}$ [ $A \cdot$ ][ ]$^x$ $x$ $b$ .

- Learning NNs is hard, often done via gradient descent

  - Topics skipped – "strongest" hardness results based on crypto [Klivans, Sherstov] ( hardness of improper learning ) .

  intersection of half-spaces .

- Analyzing gradient descent:

  (for NNs)

  - Can one analyze dynamics of gradient descent? [can view as Kernel regression for a time-varying kernel] re-phrasing .

  - Are there cases where we can reason about resulting solution? [for infinitely wide nets, kernel remains "fixed" – neural tangent kernel]

  Law of Large Numbers ( Concentration bounds ) .

# NTK REVIEW

**Theorem.** [Jacot, Gabriel, Hongler 18] [Arora, et al. 2019] A width ~ n^3 network (any number of layers) trained via GD from random initialization achieves zero training error. Moreover, the final solution is equivalent to solving a "Kernel regression" problem with a specific kernel.

$f$ .    $K(x,y):$ similarity between $x, y$ .

hypotheses of the form : $f(x) = \sum \alpha_i K(x, x_i)$

$i \in$ training samples .    $i^{th}$ training sample .

- Kernel regression ✓

- Any model training can be viewed as Kernel regression with time varying kernel

- With wide DNNs, kernel doesn't change much!

✗ (GD with tiny step size).

( Can be used to show that if width >> |training data|, then NN training via a GD converges to Kernel regression with NTK .).

# NTK EXPERIMENTS

$$K(x,y) = e^{-\|x-y\|^2 \cdot}$$

[Arora, et al. 2019]  What happens if we forget about NNs, compute closed form for NTK (determined only by number of layers, types of connections, activation function), perform kernel regression?
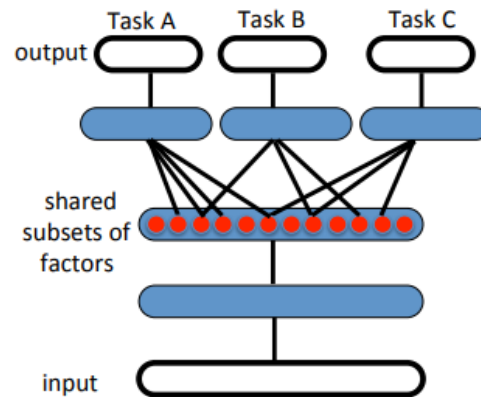
"vanilla" vrn.

| Depth | CNN-V | CNTK-V | CNN-GAP | CNTK-GAP |
|-------|-------|--------|---------|----------|
| 3 | 61.97% | 64.67% | 57.96% | 70.47% |
| 4 | 62.12% | 65.52% | 80.58% | 75.93% |
| 6 | 64.03% | 66.03% | 80.97% | 76.73% |
| 11 | 70.97% | 65.90% | 75.45% | 77.43% |
| 21 | 80.56% | 64.09% | 81.23% | 77.08% |

CIFAR-10.

# REPRESENTATION LEARNING

*(handwritten, top right)* $F_2(F_1(x)...$  $F_r(F_{r-1}...$  $F_1(x))).$  $F_1(x)$  input.

- General idea – neural networks are "hierarchical feature extractors"

- Circa 2000s – manual feature extraction (HOG, SIFT)

- NNs embed inputs -> "feature space" (alternative 'representation')

*(handwritten)* layered NN → builds iterative representation of data.  → "related" tasks  ⇝ Caption generation.  → [Bengio, et al. 2012]

Fig. 1. Illustration of representation-learning discovering explanatory factors (middle hidden layer, in red), some explaining the input (semi-supervised setting), and some explaining target for each task. Because these subsets overlap, sharing of statistical strength helps generalization.

# REPRESENTATION LEARNING

features ⟷ representation.
space embedding

- <u>What makes a good representation?</u>

  - Contrastive (for classification)

  - "Disentangled" or orthogonal

  - Sparse "explanations" for phenomena

  - Hierarchically organized, explanatory

  to robustness . . .

  in this "feature space", inputs get clustered based on class.

  → classification.

  basis in which inputs have a "sparse" representation.

- <u>Supervised vs Unsupervised</u>

  → no ta

  → Crucial if we want to (do) multiple tasks with same representation.

  animal      household object      boat.

  → Word embeddings..in NLP.

  → Graph vertex embeddings.

# REPRESENTATION LEARNING

given data ; no task ;

no labels.

■ Unsupervised learning of representation

■ Sparse coding / autoencoders (past)

■ Self-supervision (present/future)

"Meta qn": ~~Under~~ Want to "understand" data …

Formalizing : find comōn patterns / compression.

succinct representation.

Sparse Coding : Is there a "basis" for data in which all the data points are ~ sparse? (in the lin·alg. sense).

Inputs: $\qquad x_1, x_2, \ldots, x_N \in \mathbb{R}^d$.

Can you find a "basis", i.e, $v_1, v_2, \ldots, v_m \in \mathbb{R}^d$ such that

every $x_i \approx \sum_j \alpha_j^{(i)} v_j$, for some "sparse" $\alpha^{(i)}$?

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ (at most $k$ are non-zero).

\# parameters in "input rep": $dN$

\# parameters in "new" rep: $\cancel{k_1 + m}$ $\boxed{kN + md}$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ ↓

$\qquad\qquad\qquad\qquad\qquad\qquad$ coeffs.

$\qquad$ Want $\quad kN + md \ll dN$.

JPEG: based on ~ ideas.-

$\qquad$ ⟶ Layerwise unsupervised pre-training (2013..).