

15-780: Lecture 3

Aditi Raghunathan

Jan 22 2023

Recap

- We are building towards GPT4 and CLIP

- Supervised learning

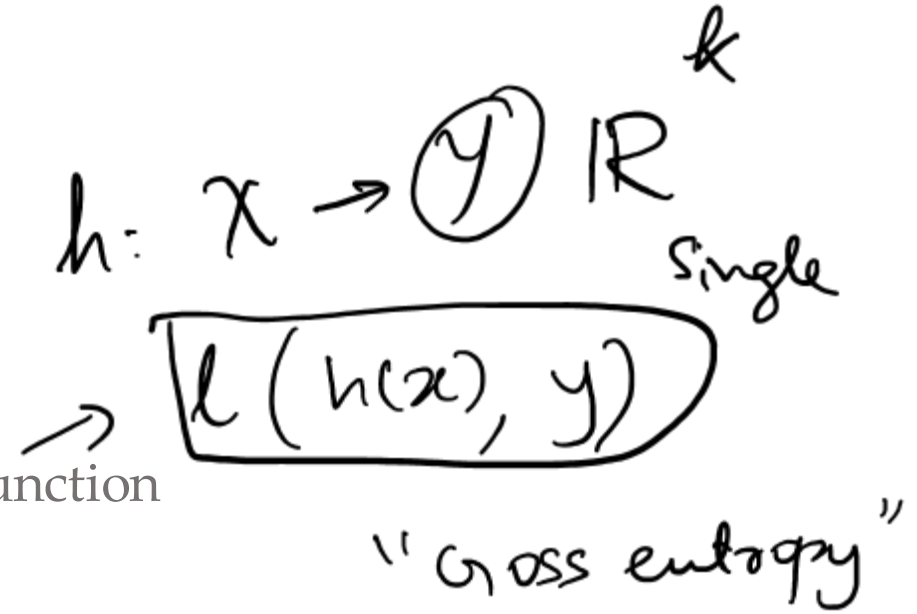
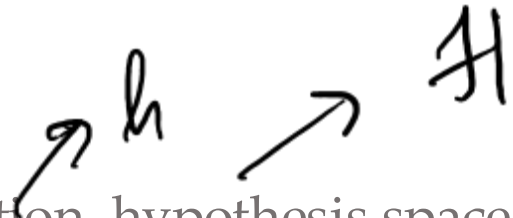
- Mapping input to targets

Notation: hypothesis function, hypothesis space, loss function

- Minimizing training loss

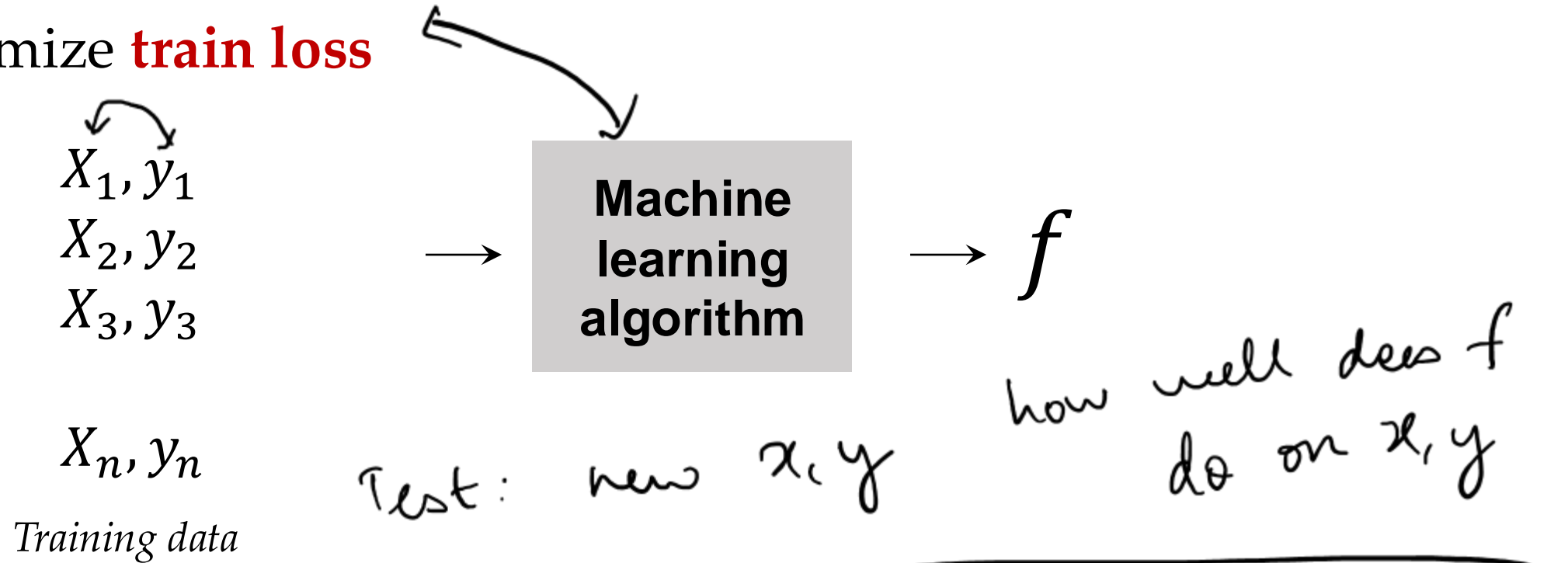
- Why does this work? Generalization (This lecture)**

- How to minimize training loss? Optimization



Minimize training loss

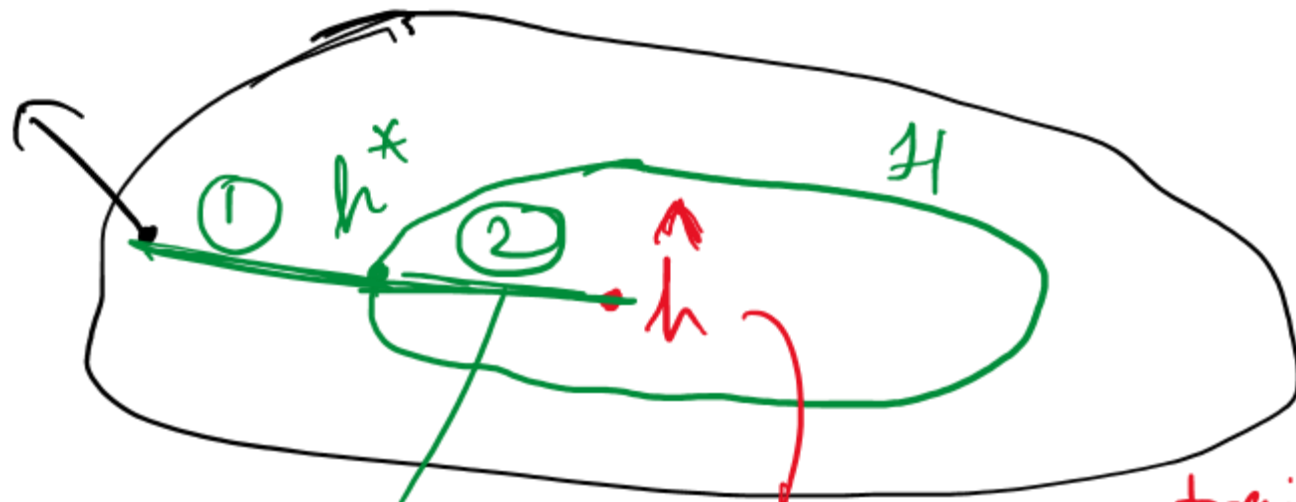
- Minimize **train loss**



We care about good performance on **unseen examples** (test set)

Intuitive picture

"best" f^*
possible



train \neq test

minimizing train loss

① $h^* \neq f^*$

② $\hat{h} \neq h^*$

\mathcal{H} : set of hypothesis function h
 $h: \mathcal{X} \rightarrow \mathbb{R}^k$ $l(h(x), y) \in \mathbb{R}$

p^* : underlying distribution over \mathcal{X}, \mathcal{Y}

Train data: $x^{(i)}, y^{(i)} \stackrel{iid}{\sim} p^*$ (n of these)

Test data: $x, y \sim p^*$

Expected risk

$$L(h) = \mathbb{E}_{x, y \sim p^*} [l(h(x), y)] \Rightarrow \text{we care about this}$$

Expected risk minimizer

$$h^* \in \underset{H}{\operatorname{argmin}} L(h)$$

"best possible hypothesis"

Empirical risk

↓ observed (training samples)

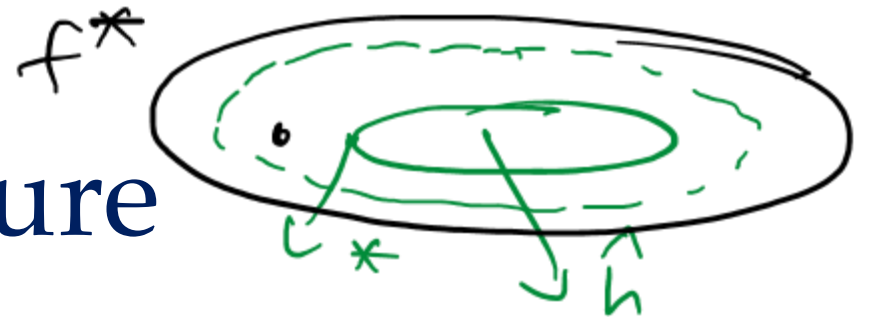
$$\hat{\Delta}(h) = \frac{1}{n} \sum_{i=1}^n \ell(h(x^{(i)}, y^{(i)}))$$

Empirical risk minimizer

$$\hat{h} \in \underset{\mathcal{H}}{\operatorname{argmin}} \hat{L}(h)$$

$$\hat{h} \neq h^*$$

Formalizing intuitive picture



- Approximation error

$$L(h^*) - L(f^*)$$

purely a fn of \mathcal{H}

- Estimation error

$$L(\hat{h}) - L(h^*)$$

"excess risk"

$$= L(\hat{h}) - L(f^*)$$

$$= \boxed{L(\hat{h}) - L(h^*)}$$

$$+ L(h^*) - L(f^*)$$

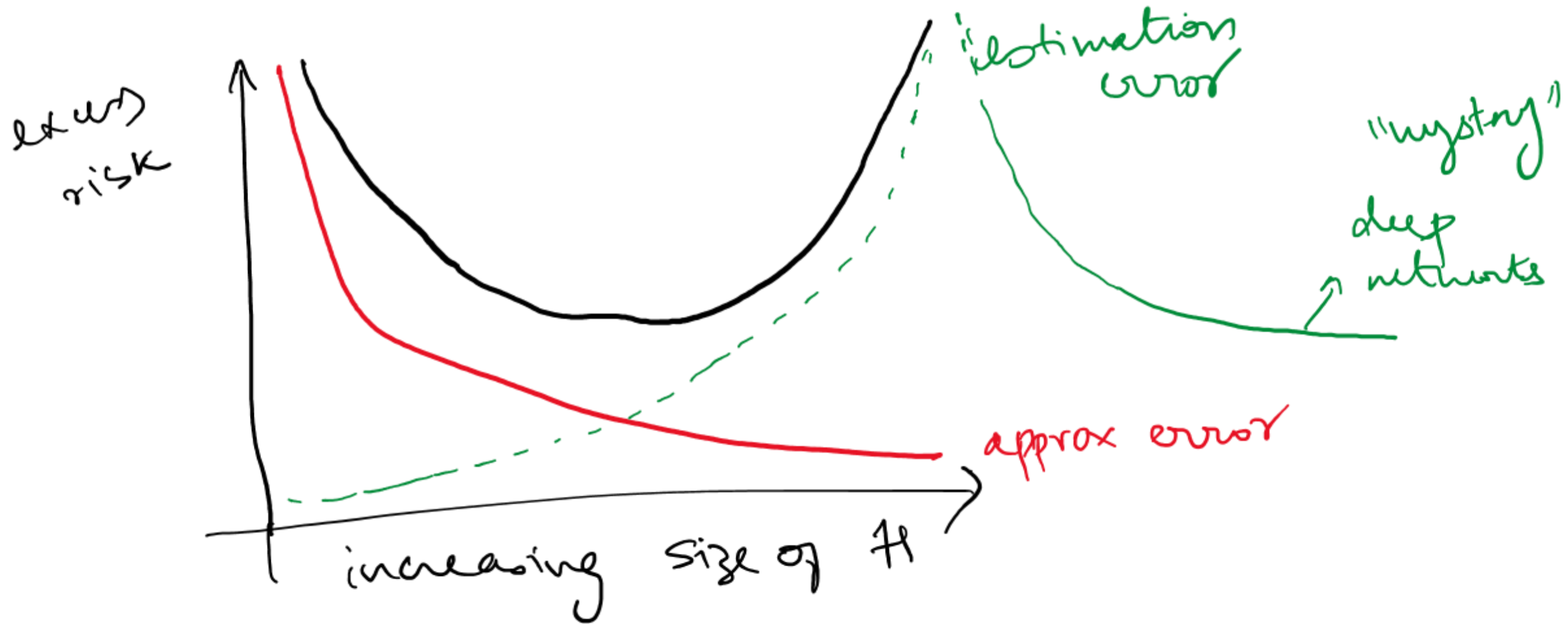
approx error



Piazza poll



Effect of hypothesis class size



Understanding estimation error

Estimation error: $L(\hat{h}) - L(h^*)$

$L(\hat{h})$: a random quantity

$$\mathbb{P} \left[\underbrace{L(\hat{h}) - L(h^*)}_{\text{estimation is high}} > \epsilon \right] < \delta$$

low probability

Simple realizable case

① $h: \mathcal{X} \rightarrow \{0, 1\}$ "deterministic"

② loss function: zero-one error $\mathbb{1}[h(x) \neq y]$

③ \mathcal{H} is finite

④ "realizable": $f^* = h^*$
 $\exists h^* \text{ s.t. } \forall x, y$

$$\Rightarrow \hat{L}(h^*) = 0 \Rightarrow \hat{L}(\hat{h}) = 0 \quad h^*(x) = y$$

we want to bound $P[L(\hat{h}) - L(h^*) > \epsilon]$

$$L(\hat{h}) - \hat{L}(\hat{h}) + \hat{L}(\hat{h}) - L(h^*)$$

\hat{h} is defined as minimizer of \hat{L}

$B: \{h \mid L(h) > \epsilon\}$ set of bad hypotheses

$$P[L(\hat{h}) - L(h^*) > \epsilon] = P[L(\hat{h}) > \epsilon]$$

\circ (realizability) $= P[\hat{h} \in B]$

$$P[\hat{h} \in B]$$

$$\hat{L}(\hat{h}) = 0$$

Step one: If hypothesis which is bad $[L(h) > \epsilon]$

what is probability it has zero train loss

$$P[\hat{L}(h) = 0]$$

$$(1 - \epsilon)^n$$

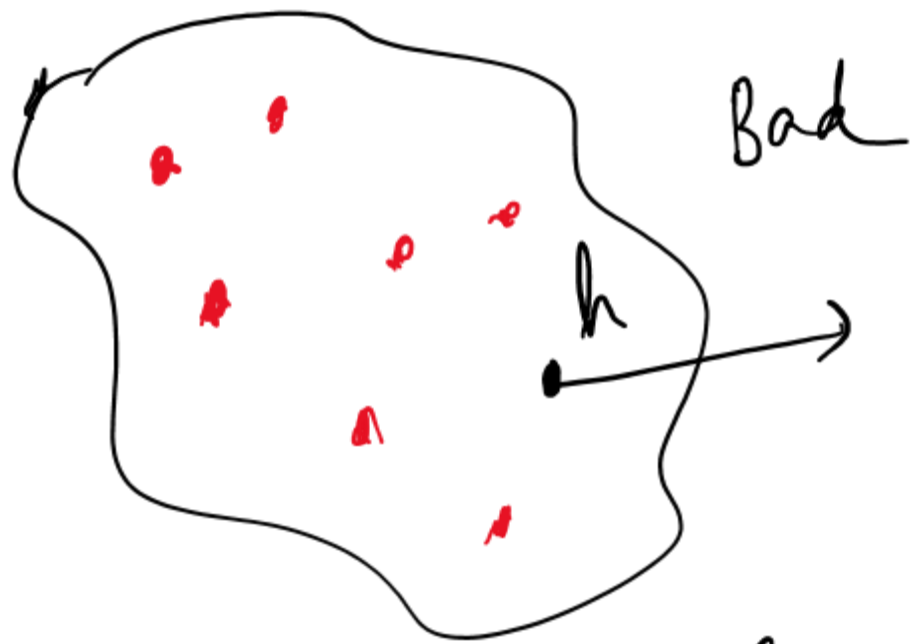
$$\leq e^{-\epsilon n}$$



ϵ mass where h is wrong

$$1 - x \leq e^{-x}$$

qs $P[\hat{h} \in B] \leq (1 - \epsilon)^n ?$



"appears good" $(1-\epsilon)^n$

h : any classifier that has zero train loss
or appears good $P(\hat{h} \in B)$?

union bound: $P[\exists h \in B; \hat{L}(h) = 0] \leq \sum_{h \in B} P[\hat{L}(h) = 0] \leq |B| e^{-\epsilon n} \leq |H| e^{-\epsilon n}$

$$P(\hat{h} \in B)$$

$$\leq |H| e^{-\epsilon n} \delta$$

union bound

appears good on n samples despite being bad

estimation error $> \epsilon$

$$\text{w.p. } 1 - \delta \quad L(\hat{h}) \leq \frac{\log(H) + \log(1/\delta)}{n}$$

excess risk
= estimation error

as n increases, est error decreases
as (H) increases, est error increases

General recipe

- **Convergence:** for fixed h , $L(h)$ is close to $\hat{L}(h)$ test loss train loss as $n \uparrow$, gap goes down
- **Uniform convergence:** convergence holds for all hypothesis simultaneously ↓
makes it harder as H expands
- Why uniform convergence?

Takeaways

- Approximation error: decreases with increase in H
- Estimation error: more nuanced, depends on H
 - Very large H leads to high estimation error
- How to keep H small?

Regularization

- Linear classifiers: dimensionality, norm

$$h(x) = \phi^T x \quad \phi \in \mathbb{R}^d \quad \|\phi\| \text{ small}$$

- Regularized objective

objective: $\text{Train loss} + \lambda \|\phi\|_2$

Any questions?

