# 9. Bias-Variance Tradeoff

Error Decomposition aka Bias-Variance Tradeoff    §5.2 [SS]

Let $h_S \in ERM_{\mathcal{H}}(S)$

Total Error (true risk) is

$$L_D(h_S) = \underbrace{L_D(h_S) - \min_{h \in \mathcal{H}} L_D(h)}_{\text{estimation error}} + \underbrace{\min_{h \in \mathcal{H}} L_D(h)}_{\text{approximation error}}$$

$$\overset{\varepsilon_{est.}}{} \qquad \overset{\varepsilon_{approx.}}{}$$

or

$$\underbrace{L_D(h_S) - \underbrace{\min_{\text{all } h} L_D(h)}_{\text{Bayes risk, } \varepsilon_{Bayes}} = \varepsilon_{est} + \widetilde{\varepsilon}_{approx}}_{\text{"excess risk"}}$$

$$\widetilde{\varepsilon}_{approx} = \varepsilon_{approx} - \varepsilon_{Bayes}$$

$\varepsilon_{est}$: error due to using $\hat{L}_S$ instead of $L_D$. This is what our agnostic PAC bounds cover.

   aka generalization.    Results like $m = O\left( \frac{\log(|\mathcal{H}|/\delta)}{\varepsilon^2} \right) \cdots$

Smaller $|\mathcal{H}|$ is "good": better sample complexity, generalization, lower variance

   in fact, soon we'll see metrics to deal with $|\mathcal{H}| = \infty$
   The intuition is low complexity $\mathcal{H}$ is good

$\varepsilon_{approx}$: we've ignored so far. This is more classical, or unknowable a priori.   Ex papayas



we're looking to approximate $h_{true}$ by $h \in \mathcal{H}$

Ex    $\mathcal{H} = \{$ all polynomials of degree $\leq 100\}$

   $h_{true}$ is an arbitrary function (maybe even continuous)
   $\varepsilon_{approx}$ unlikely 0

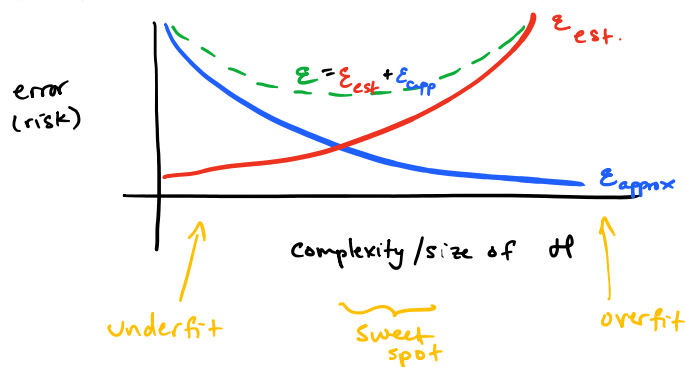$\mathcal{E}_{Bayes}$ accounts for inherent noise in labels

eg, (pure) PAC: $y = f(x)$ so $\min_h L_D(h) = L_D(f) = 0$

eg, agnostic PAC: $y \sim D(y|x)$

eg $y = f(x) + z$, $z \sim N(0, \sigma^2)$
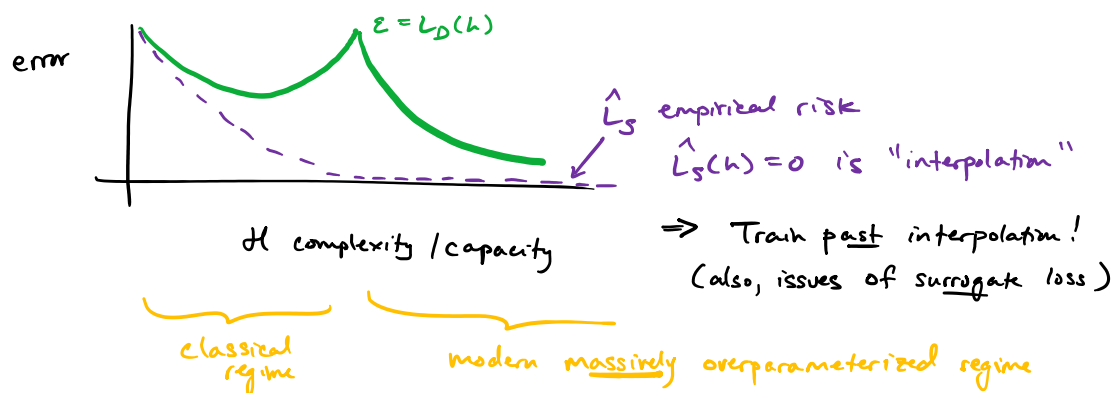
not our focus

(textbook) bias-variance tradeoff



error (risk)

$\mathcal{E}_{est.}$

$\mathcal{E} = \mathcal{E}_{est} + \mathcal{E}_{app}$

$\mathcal{E}_{approx}$

Complexity/size of $\mathcal{H}$

underfit

sweet spot

overfit

Point: both terms matter, so $|\mathcal{H}|$ small not always good

$\mathcal{E}_{approx}$ hard to predict a priori

## Double-descent

Belkin, Hsu, Ma, Mandal PNAS '19



error

$\mathcal{E} = L_D(h)$

$\hat{L}_S$ empirical risk

$\hat{L}_S(h) = 0$ is "interpolation"

$\Rightarrow$ Train past interpolation!

(also, issues of surrogate loss)

$\mathcal{H}$ complexity/capacity

classical regime

modern massively overparameterized regime

Classical: "A model with 0 training error is overfit to the training data and will typically generalize poorly"

(p 221, Hastie, Tibsirani, Friedman "The Elements of Statistical Learning" 2001)

That <u>can</u> be true, but need not be

- Deep Learning (2013+) empirical results strongly show best results obtained on massively overparameterized ($\hat{L}_S(h) = 0$)

neural networks

(caveat: how you train matters, e.g., $\arg\min_h \hat{\mathcal{L}}_S(h)$ is
a big set, and not all ERM sol'n generalize...
but the ones we find via SGD work well)

- Double descent can show up
  Though I'd disagree that it always does

We don't fully understand deep learning yet
many ideas in lit.
That paper suggests
- our notions of complexity not a good measure

- smoothness of $h$ might be better
  (already sometimes exploited)

  " $H$ larger might allow for smoother functions
  (exactly at interpolation threshold is likely not smooth)

- or, w/ many sol'n to choose from, a least-squares
  (or other crit.) works well