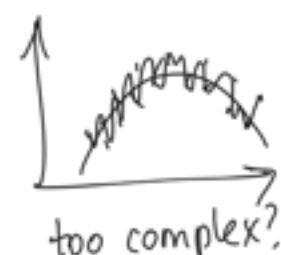


Model Selection

Today: Non-Bayesian methods.



core challenge in ML: how to generalize?
→ want to predict well in the future!

mostly focus on how to identify if a model will generalize
(rest of course: ways to generalize)

Examples:

HW: sunspots vs Republicans in Senate

Simple case:

- suppose we have $N=8$ datapoints

- data: x_n : 2000 dims, all binary

all $x_{nd} \sim p$, $p = 0.5$

$y_n = x_{n1}$ (perfect predictor of y given x)

- What is the probability that $y_n = x_{n2}$ for all n in $1..N$
($N=8$): $(\frac{1}{2})^N = (\frac{1}{256})$

- What is the probability that there is at least one spurious perfect correlation? ($d \neq 1$)

$$\Pr(\text{at least one spurious}) = 1 - \Pr(\text{no spurious at } d) = 1 - \left(\frac{255}{256}\right)^{1999} \approx 1$$

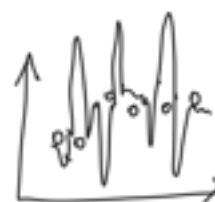
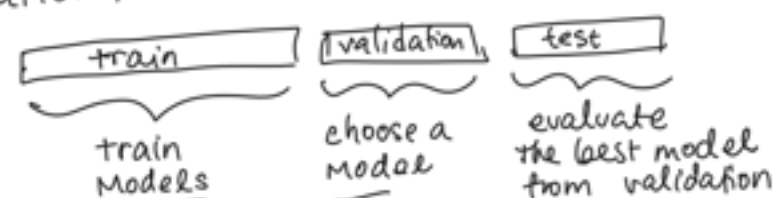
What can we do? Today: statistical techniques, but

1) always inspect for plausibility

2) this will not discover issues due to confounds
stat methods

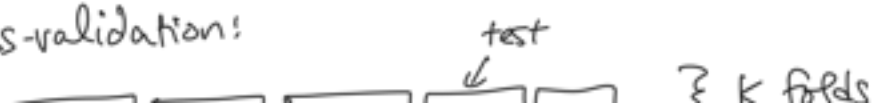
(assume test data is like train data)

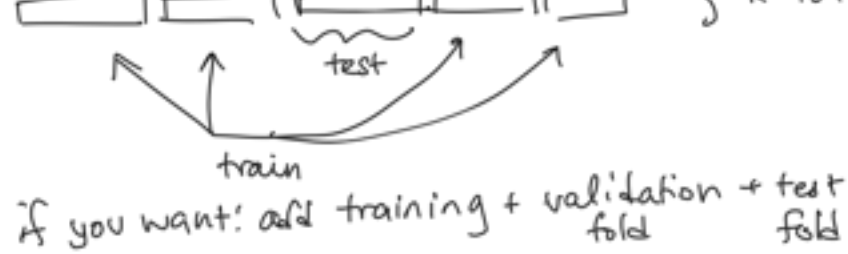
Validation / Cross-Validation



What if we want confidence intervals?
[And not enough data]

↳ cross-validation:





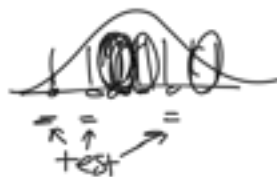
(not indep, but can give sense of variance)

→ Bootstrap: where we sample with replacement N samples from our data; train; test on rest

stat story: $X, Y \sim \mathcal{D}$

if we train: $\hat{y} = f(x)$, err } this is the frequentist question
 $\mathbb{E}_{\mathcal{D}}[\text{err}] \quad \mathbb{E}_{\mathcal{D}}[\text{err}^2]$

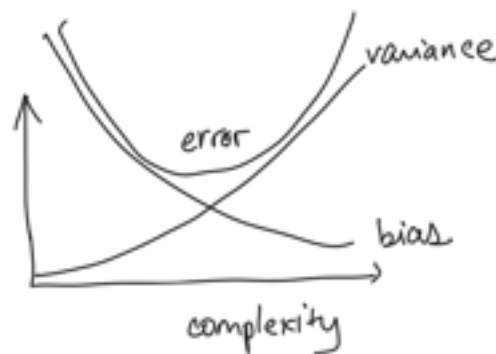
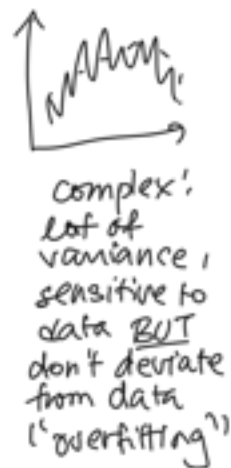
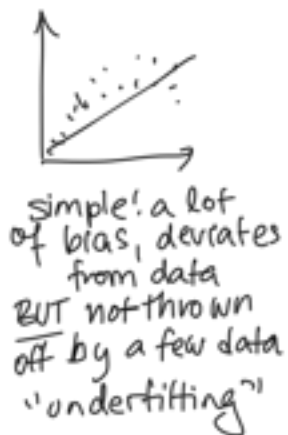
bootstrap: ~~data~~ $X, Y \sim$ subset of $[X, Y]$



⇒ All of these help measure generalization (empirical approach)

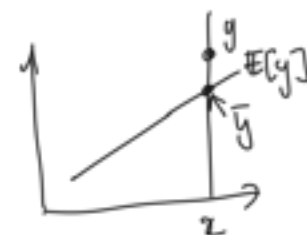
Deeper look at WHY models generalize

bias-variance trade-off



Let's formalize! We care about $\mathbb{E}_{\mathcal{D}}[(y - \hat{y})^2] = \mathbb{E}_{\mathcal{D}}[(y - f_{\mathcal{D}}(x))^2]$
 for some x , assume some true process to create y 's } squared error
 \uparrow predictor depends on data \mathcal{D}

real life $[y$ is a random variable with mean $\bar{y} = \mathbb{E}[y]$



$$\begin{aligned} \mathbb{E}[(y - f_{\mathcal{D}}(x))^2] &= \mathbb{E}[(y - \bar{y} + \bar{y} - f_{\mathcal{D}}(x))^2] \\ &= \underbrace{\mathbb{E}[(y - \bar{y})^2]}_{\text{noise inherent in the real situation}} + \underbrace{\mathbb{E}[(\bar{y} - f_{\mathcal{D}}(x))^2]}_{\text{error from our modeling}} + 2\underbrace{\mathbb{E}[(y - \bar{y})(\bar{y} - f_{\mathcal{D}}(x))]}_{\text{0 mem indep.}} \end{aligned}$$

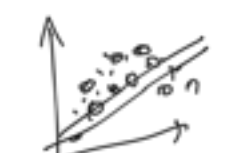
⇓ expand:

lets consider all the different $f_{\mathcal{D}}(x)$ we could

learn w/ diff datasets \mathcal{D} .

$$\mathbb{E}_{\mathcal{D}} [\bar{y} - f_{\mathcal{D}}(x)]^2 = \underbrace{\mathbb{E}_{\mathcal{D}} [\bar{y} - \bar{f}(x)]^2}_{\substack{\text{real world} \\ \text{expectation,} \\ \text{constant}}} + \underbrace{\mathbb{E}_{\mathcal{D}} [\bar{f}(x) - f_{\mathcal{D}}(x)]^2}_{\substack{\text{depends on} \\ \mathcal{D}, \text{ define} \\ \bar{f}(x) = \mathbb{E}_{\mathcal{D}} [f_{\mathcal{D}}(x)]}} + \underbrace{2\mathbb{E}_{\mathcal{D}} [\bar{y} - \bar{f}(x)](f_{\mathcal{D}}(x) - \bar{f}(x))}_{\substack{\text{indep} \\ \text{mean 0}}} + \underbrace{\mathbb{E}_{\mathcal{D}} [f_{\mathcal{D}}(x) - \bar{f}(x)]^2}_{\substack{\text{variance,} \\ \text{what you} \\ \text{get depends} \\ \text{on } \mathcal{D}}}$$

bias²
mean of $f_{\mathcal{D}}(x)$
doesn't match
real mean!



$$\text{error} = \text{noise}^2 + \text{bias}^2 + \text{variance}$$

control

Note: $N \rightarrow \infty$ (or large data), var \downarrow

How can we manage bias, variance trade-off?

→ Regularization: penalize model complexity

$$\text{e.g. } \mathcal{L}_{\mathcal{D}}(w) \rightarrow \min_w \mathcal{L}_{\mathcal{D}}(w) + \lambda R(w)$$

$$R(w): \|w\|_2^2 \quad \text{Ridge Reg. / L2}$$

$$\|w\|_1 \quad \text{L1/lasso}$$



convex :)

convex,
non-diff.

choose λ via cross-validation.

→ Ensembles: use a committee to make a decision
(vote/avg)

Idea: bias doesn't increase/change

variance goes down

→ (~~BAGGING~~) ⇒ sample datasets

→ RF / extra random forests



→ Adaboost / Boosting

sequentially learn predictors that are indep of prev predictors

Concept Check:

Test Set & Train set for some data

Model A

Model B

test: 0.70

accuracy

test: 0.60
train: 0.95 train: 0.72

- 1) Will adding more data help for model A? model B?
- 2) Should we have a validation split in our train?
- 3) Which model will have stable or similar results if we try on another dataset to train.