

What is learnable?

Readings for today

- Fulop, S. A., & Chater, N. (2013). Learnability theory. *Wiley Interdisciplinary Reviews: Cognitive Science*, 4(3), 299-306.

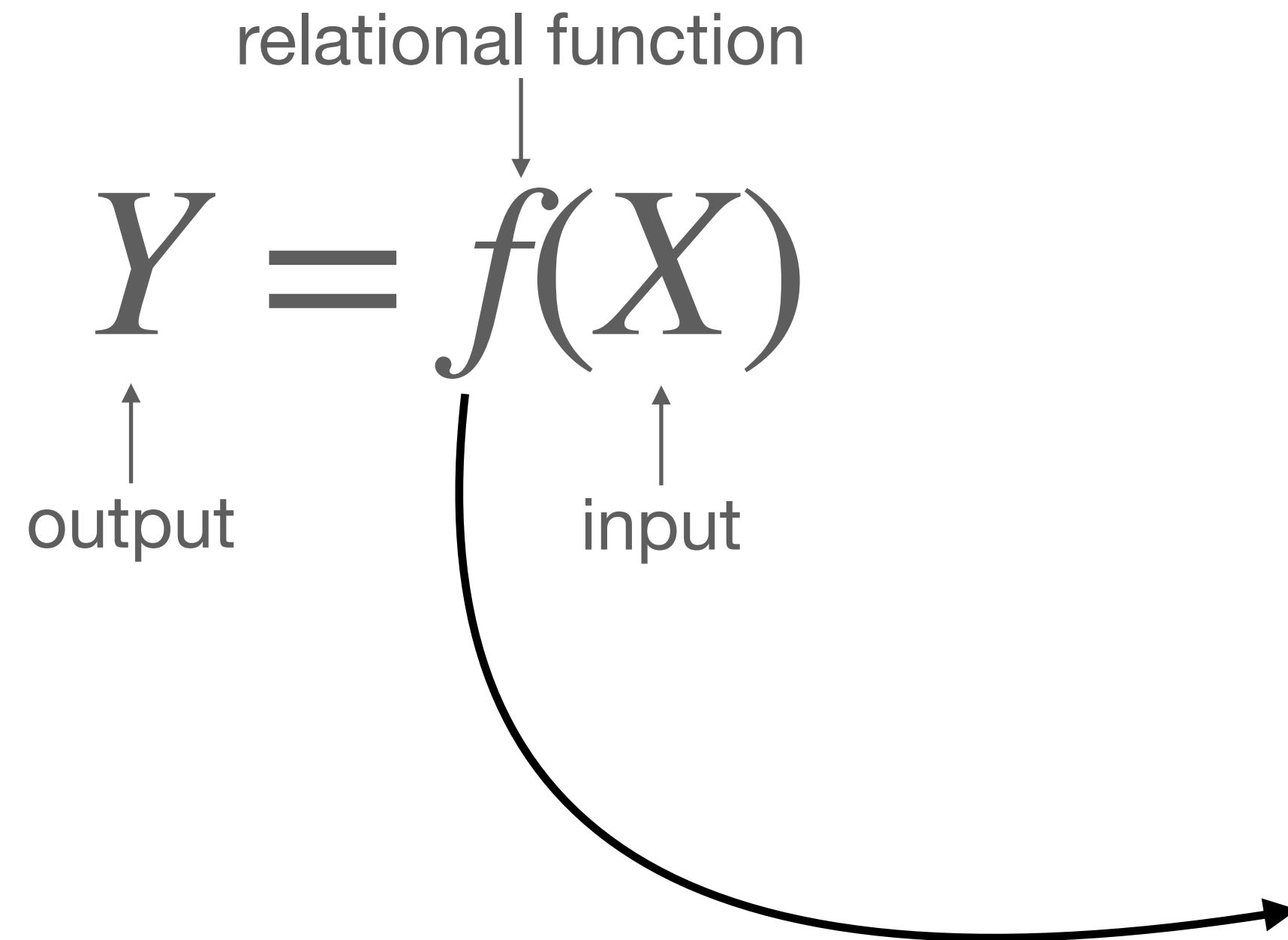
Topics

1. The two parts of a statistical model
2. Understanding what is learnable

The two parts of a statistical model

Fundamental form of a statistic

Fundamental form of a statistic



Sample

$$S = \{(X_i, Y_i)\}_{i=1}^n$$

- h : hypothesis \rightarrow prediction based on X
- g : learner \rightarrow system that finds the best h for a particular S

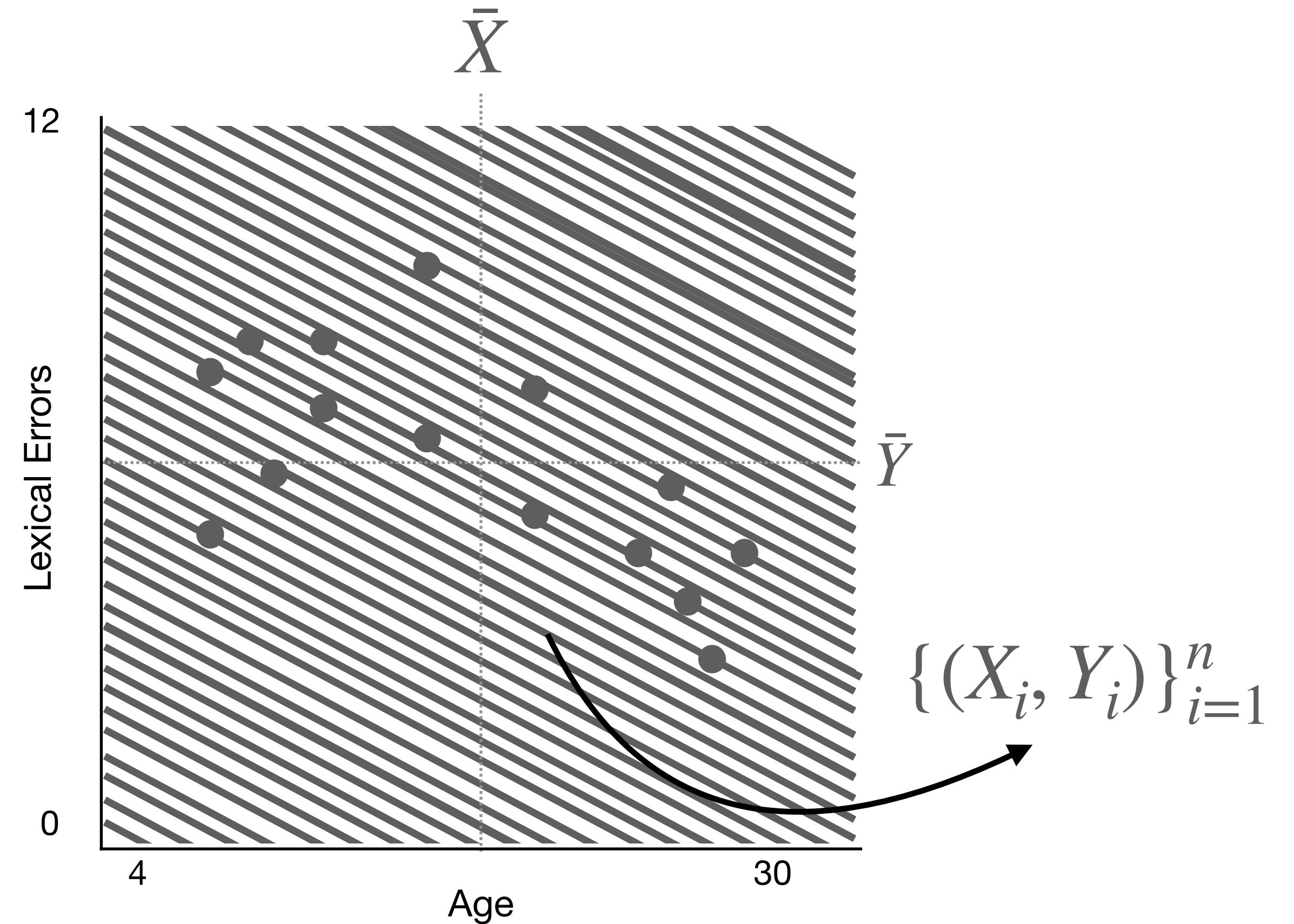
The learning problem: $h \text{ \& } g \rightarrow f$

Hypothesis

$$h : Y = \hat{\beta}_1 X + \hat{\beta}_0$$

Learner

$$g : \beta = (X'X)^{-1}X'Y$$



Empirical Risk Minimization

Risk: how well a hypothesis does in practice with new (test) data.

$$R(h) = \underbrace{\ell(h(X), Y)}_{\text{loss function}}$$

Goal of learner

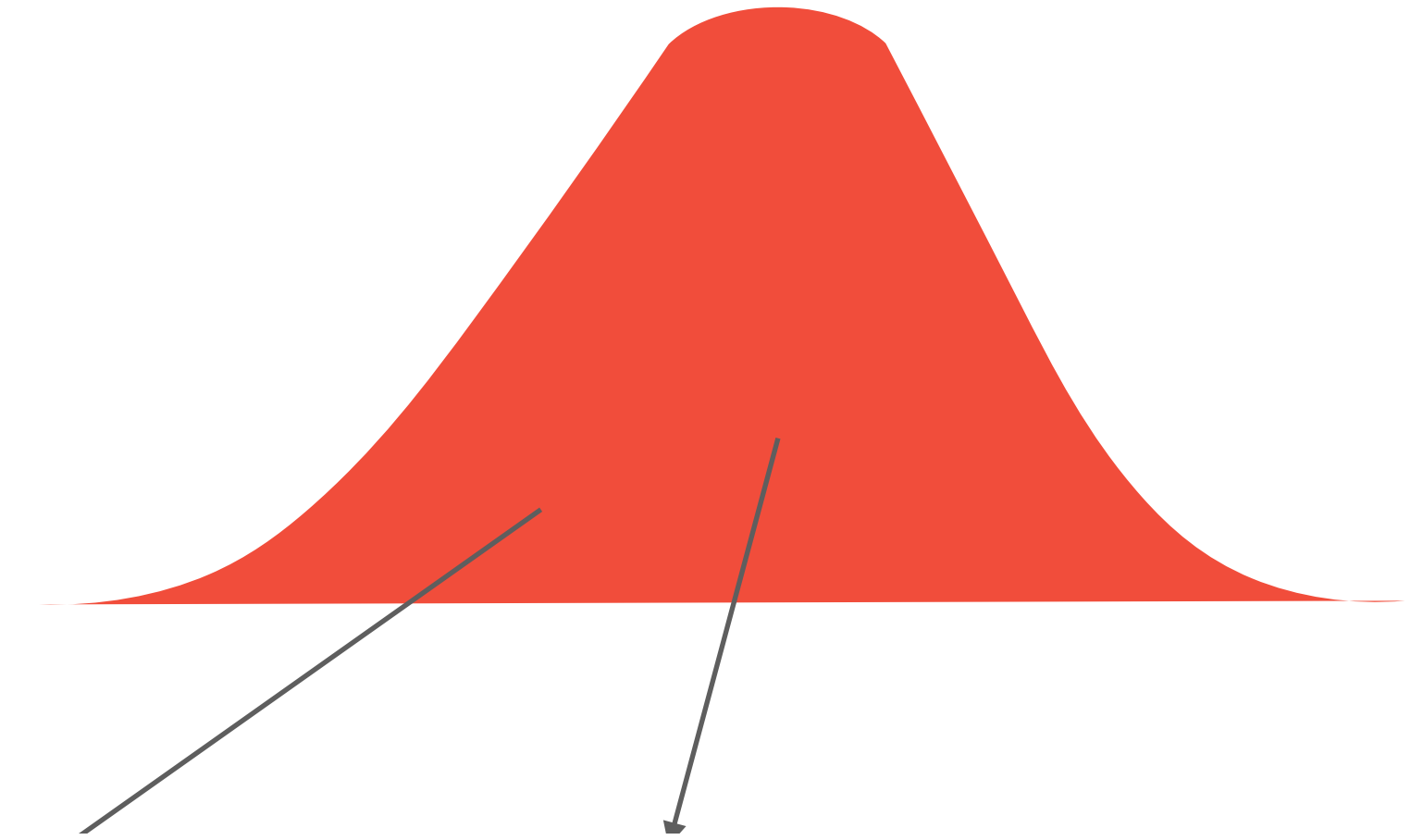
$$g = \arg \min_{\underbrace{h \in H}} R(h)$$

$$\ell(h(X), Y) = \begin{cases} 1, & h(X) \neq Y \\ 0, & h(X) = Y \end{cases}$$

Empirical Risk Minimization

Expected Risk

$$E_{\text{risk}}(h, n, P) = \int_{\underbrace{(\mathbf{X}, \mathbf{Y})_n}_{\text{train}}} \underbrace{R(h)}_{\text{risk}} \underbrace{dP_{(X, Y)_n}}_{\text{distribution}}$$



Assumption: Both the training and test data come from the same distribution.

Understanding what is learnable

What is learnable?

Learnability theory: Can the true h be learned?

Determine whether the true h is able to be learned given a particular data scenario (e.g., particular $P_{(\mathbf{X}, \mathbf{Y})}$)?

Computational theory: Can the solution for the true h be computed?

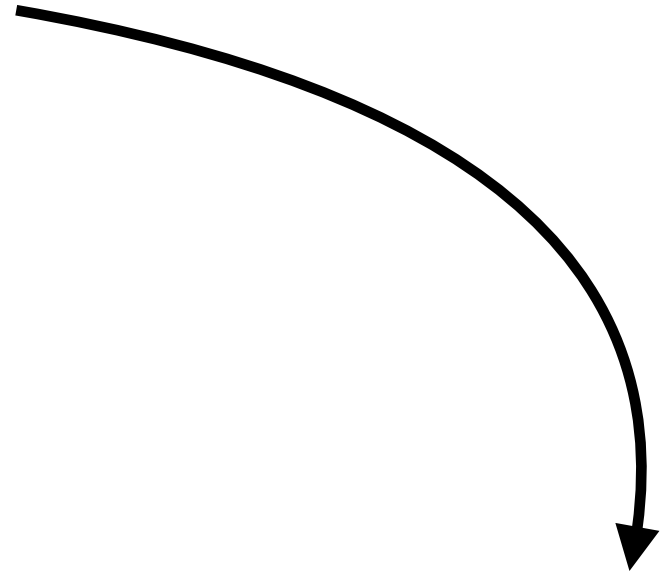
P : Problem can be *computed* in polynomial time.

→ We can only work with
“ P hard” problems

NP : Problem requires non-deterministic polynomial time.

Bayesian Learning

Bayes Rule:

$$P(H_i | \underbrace{\{(X, Y)\}_{i=1}^n}_D) = \frac{P(D | H_i) P(H_i)}{P(D)}$$


H_i : individual hypothesis

Naive prior:

$$P(H_i) = U(-\infty, +\infty)$$

$$P(H_i | D) = P(D | H_i)$$

Bernstein-von Mises Theorem

- If:
- (X, Y) are independent and identically distributed
 - limited to a finite space of outcomes.

Assuming: • The true h is learnable at all.

Then: • The posterior probability converges to the true h among almost all samples.

 e.g., testing & training sets

Model selection

$$h_{true} \rightarrow h_{best}$$

Generalization: Given a fixed **training data set**, find the model that *best* predicts future **unseen (test) data set**.

$$E_{\text{risk}}(h, n, P) = \underbrace{\int_{(\mathbf{X}, \mathbf{Y})_n}}_{\text{train}} \underbrace{\int_{(\mathbf{X}, \mathbf{Y})}}_{\text{test}} \underbrace{\ell(h(\mathbf{X}), \mathbf{Y})}_{\text{loss function}} \underbrace{dP_{\mathbf{X}, \mathbf{Y}}}_{\text{distribution}} \underbrace{dP_{(\mathbf{X}, \mathbf{Y})_n}}_{\text{distribution}}$$

Model selection

Bayesian evidence:

Relative evidence in favor of H_m over H_j

$$\frac{P(H_m | D)}{P(H_j | D)} = \frac{P(H_m)}{P(H_j)} \times \frac{P(D | H_m)}{P(D | H_j)}$$

Naive prior:

$$P(H_i) = U(-\infty, +\infty)$$

Neyman-Pearson Lemma:

When probabilities are known, the optimal decision rule for selecting one H over another is a likelihood ratio test.

Probably Approximately Correct (PAC) Learning

Q: How do you get “good enough” learning so as to be useful?

PAC learning requires a learner to:

1. Approximate the true h
2. Be computationally feasible (P problem)

Approximately: A hypothesis $h \in H$ is approximately correct if its error over the training data $P_{(X,Y)}$ is bounded by some ϵ , with $0 \leq \epsilon \leq \frac{1}{2}$.

Probably: The h is probably approximately correct at a generalization error rate of δ , if its prediction accuracy is $1 - \delta$, with $0 \leq \delta \leq \frac{1}{2}$.

Probably Approximately Correct (PAC) Learning

Bounding the sample size:

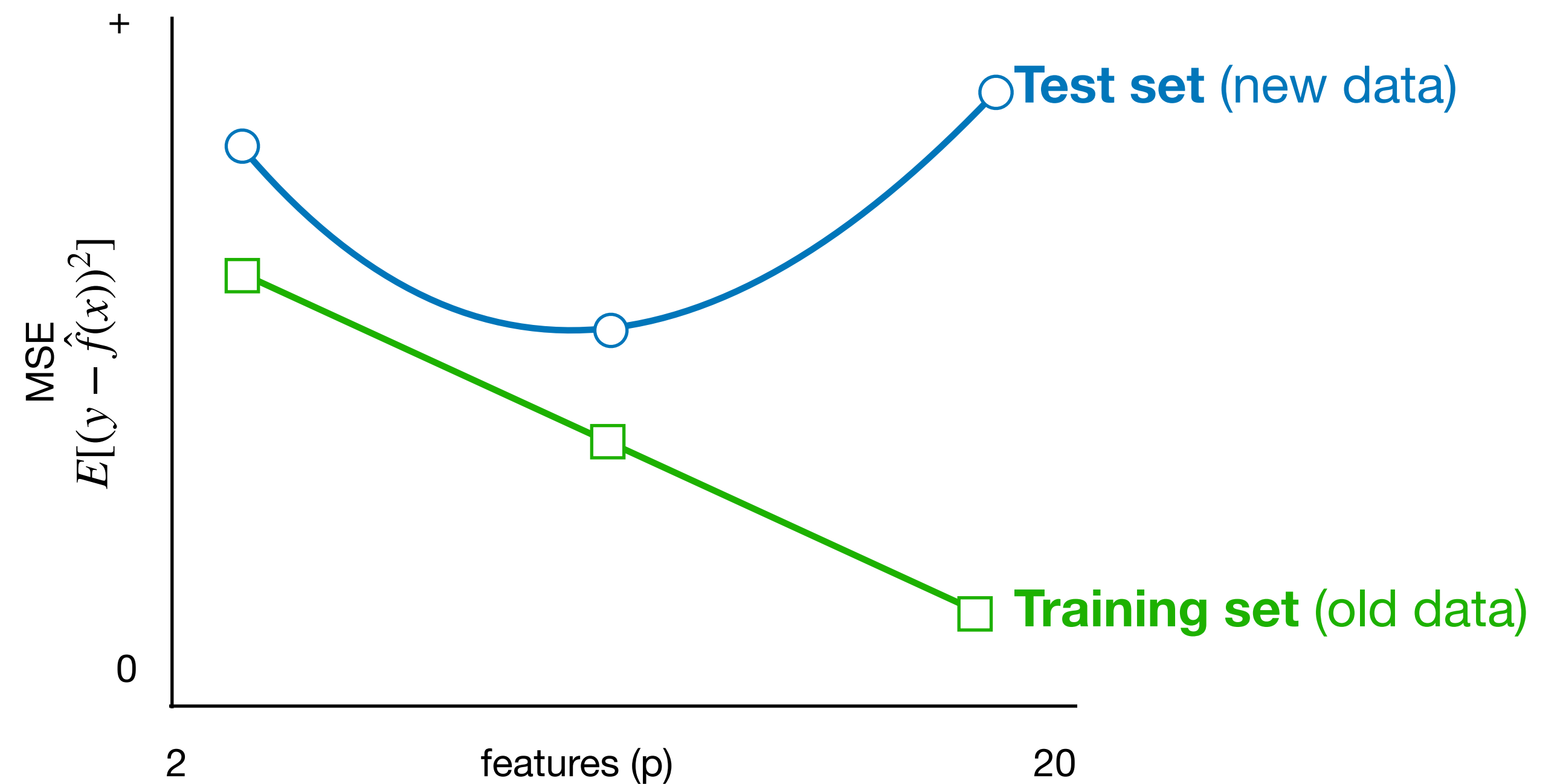
$$n \geq \frac{1}{\epsilon} \left(\ln |H| + \ln \frac{1}{\delta} \right)$$

• as $\downarrow \epsilon$ & $\downarrow \delta$, or $\uparrow H$, then $\uparrow n$

e.g., complexity of model

The bias-variance tradeoff:

- $\delta \rightarrow$ Test error
- $\epsilon \rightarrow$ Training error



Take home message

- A statistical model consists of two parts: 1) a hypothesis, h , that predicts an output from an input and 2) a learner, g , that finds the best h for a given sample.
- Before running any analysis, you first must confirm whether your h (aka- model) is learnable and computable given the expected data set it will be applied to.