Lecture 4
Bias-Variance Tradeoff
&
Intro to Probabilistic Inference

How to quantify expected *out-of-sample* error?

$$y = f(x) + \epsilon \quad \text{with} \quad \epsilon \in \mathcal{N}(0, \sigma_\epsilon^2)$$

$$\ell(y, h(X; \boldsymbol{\theta})) = \sum_i \left( y_i - h(x_i; \boldsymbol{\theta}) \right)^2$$

Given a fixed training set $D$, the empirically optimal hypothesis is parametrised by

$$\hat{\boldsymbol{\theta}}_D = \arg \min_{\boldsymbol{\theta}} \ell(y, h(X; \boldsymbol{\theta}))$$

On average, how well does our empirically optimal hypothesis $h_S^* \equiv h\left(x, \hat{\boldsymbol{\theta}}_D\right)$ *make prediction*?

In other words, we want to quantify how well our empirically optimal hypothesis from finite samples predict the value of the unknown god-given function $f$.

We'll perform a bias-variance decomposition to do quantify this.

Let's calculate the out-of-sample error, averaged over multiple data samples $D$ and noise realisation $\epsilon$

$$\mathbb{E}_{D,\epsilon}\left[\ell\left(\boldsymbol{y}, h\left(\boldsymbol{X}; \hat{\boldsymbol{\theta}}_D\right)\right)\right] = \mathbb{E}_{D,\epsilon}\left[\sum_i \left(y_i - h\left(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D\right)\right)^2\right]$$

Let's calculate the out-of-sample error, averaged over multiple data samples $D$ and noise realisation $\epsilon$

$$\mathbb{E}_{D,\epsilon}\left[\ell\left(\boldsymbol{y}, h\left(\boldsymbol{X};\hat{\boldsymbol{\theta}}_D\right)\right)\right] = \mathbb{E}_{D,\epsilon}\left[\sum_i \left(y_i - h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right)^2\right]$$

$$= \mathbb{E}_{D,\epsilon}\left[\sum_i \left(y_i - f\left(\boldsymbol{x}_i\right) + f\left(\boldsymbol{x}_i\right) - h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right)^2\right]$$

Let's calculate the out-of-sample error, averaged over multiple data samples $D$ and noise realisation $\epsilon$

$$\mathbb{E}_{D,\epsilon}\left[\ell\left(\boldsymbol{y}, h\left(\boldsymbol{X}; \hat{\boldsymbol{\theta}}_D\right)\right)\right] = \mathbb{E}_{D,\epsilon}\left[\sum_i \left(y_i - h\left(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D\right)\right)^2\right]$$

$$= \mathbb{E}_{D,\epsilon}\left[\sum_i \left(y_i - f\left(\boldsymbol{x}_i\right) + f\left(\boldsymbol{x}_i\right) - h\left(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D\right)\right)^2\right]$$

$$= \sum_i \mathbb{E}_{\epsilon}\left[\left(y_i - f\left(\boldsymbol{x}_i\right)\right)^2\right] + \mathbb{E}_{D,\epsilon}\left[\left(f\left(\boldsymbol{x}_i\right) - h\left(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D\right)\right)^2\right]$$

$$+ 2\mathbb{E}_{\epsilon}\left[y_i - f\left(\boldsymbol{x}_i\right)\right] \mathbb{E}_D\left[f\left(\boldsymbol{x}_i\right) - h\left(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D\right)\right]$$

Let's calculate the out-of-sample error, averaged over multiple data samples $D$ and noise realisation $\epsilon$

$$\mathbb{E}_{D,\epsilon}\left[\ell\left(\boldsymbol{y}, h\left(\boldsymbol{X};\hat{\boldsymbol{\theta}}_D\right)\right)\right] = \mathbb{E}_{D,\epsilon}\left[\sum_i\left(y_i - h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right)^2\right]$$

$$= \mathbb{E}_{D,\epsilon}\left[\sum_i\left(y_i - f\left(\boldsymbol{x}_i\right) + f\left(\boldsymbol{x}_i\right) - h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right)^2\right]$$

$$= \sum_i \mathbb{E}_\epsilon\left[\left(y_i - f\left(\boldsymbol{x}_i\right)\right)^2\right] + \mathbb{E}_{D,\epsilon}\left[\left(f\left(\boldsymbol{x}_i\right) - h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right)^2\right]$$

$$+ 2\mathbb{E}_\epsilon\left[y_i - f\left(\boldsymbol{x}_i\right)\right]\mathbb{E}_D\left[f\left(\boldsymbol{x}_i\right) - h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right]$$

$$= \sum_i \sigma_\epsilon^2 + \sum_i \mathbb{E}_D\left[\left(f\left(\boldsymbol{x}_i\right) - h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right)^2\right]$$

*measurement error*          *finite-size sampling error*
*(noise)*                    *(noise independent)*

Let's look at the last term

*finite-size sampling error*

$$\sum_i \mathbb{E}_D\left[\left(f\left(\boldsymbol{x}_i\right) - h\left(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D\right)\right)^2\right] = \sum_i \mathbb{E}_D\left[\left\{f\left(\boldsymbol{x}_i\right) - \mathbb{E}_D\left[h\left(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D\right)\right] + \mathbb{E}_D\left[h\left(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D\right)\right] - h\left(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D\right)\right\}^2\right]$$

Let's look at the last term

*finite-size sampling error*

$$\sum_i \mathbb{E}_D \left[ \left( f(\boldsymbol{x}_i) - h(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D) \right)^2 \right] = \sum_i \mathbb{E}_D \left[ \left\{ f(\boldsymbol{x}_i) - \mathbb{E}_D \left[ h(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D) \right] + \mathbb{E}_D \left[ h(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D) \right] - h(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D) \right\}^2 \right]$$

$$= \sum_i \mathbb{E}_D \left[ \left\{ f(\boldsymbol{x}_i) - \mathbb{E}_D \left[ h(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D) \right] \right\}^2 \right] + \mathbb{E}_D \left[ \left\{ h(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D) - \mathbb{E}_D \left[ h(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D) \right] \right\}^2 \right]$$

$$+ 2\mathbb{E}_D \left[ \left\{ f(\boldsymbol{x}_i) - \mathbb{E}_D \left[ h(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D) \right] \right\} \left\{ h(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D) - \mathbb{E}_D \left[ h(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D) \right] \right\} \right]$$

# Let's look at the last term

*finite-size sampling error*

$$\sum_i \mathbb{E}_D\left[\left(f(\boldsymbol{x}_i) - h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right)^2\right] = \sum_i \mathbb{E}_D\left[\left\{f(\boldsymbol{x}_i) - \mathbb{E}_D\left[h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right] + \mathbb{E}_D\left[h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right] - h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right\}^2\right]$$

$$= \sum_i \mathbb{E}_D\left[\left\{f(\boldsymbol{x}_i) - \mathbb{E}_D\left[h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right]\right\}^2\right] + \mathbb{E}_D\left[\left\{h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right) - \mathbb{E}_D\left[h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right]\right\}^2\right]$$

$$+ 2\mathbb{E}_D\left[\left\{f(\boldsymbol{x}_i) - \mathbb{E}_D\left[h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right]\right\}\left\{h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right) - \mathbb{E}_D\left[h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right]\right\}\right]$$

$$= \sum_i \left(f(\boldsymbol{x}_i) - \mathbb{E}_D\left[h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right]\right)^2 + \sum_i \mathbb{E}_D\left[\left\{h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right) - \mathbb{E}_D\left[h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right]\right\}^2\right]$$

Bias$^2$                Variance

How to quantify expected *out-of-sample* error? Putting them altogether

$$
\mathbb{E}_{D,\epsilon}\left[\ell\left(\boldsymbol{y}, h\left(\boldsymbol{X}; \hat{\boldsymbol{\theta}}_D\right)\right)\right] = \sum_i \left(f(\boldsymbol{x}_i) - \mathbb{E}_D\left[h\left(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D\right)\right]\right)^2 \qquad \text{Bias}^2
$$

$$
+
$$

$$
+ \sum_i \mathbb{E}_D\left[\left\{h\left(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D\right) - \mathbb{E}_D\left[h\left(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D\right)\right]\right\}^2\right] \qquad \text{Variance}
$$

$$
+
$$

$$
+ \sum_i \sigma_\epsilon^2 \qquad \text{Measurement error ( Noise )}
$$

# How to quantify expected *out-of-sample* error? Putting them altogether

$$\mathbb{E}_{D,\epsilon}\left[\ell\left(\boldsymbol{y}, h\left(\boldsymbol{X}; \hat{\boldsymbol{\theta}}_D\right)\right)\right] = \sum_i \left(f(\boldsymbol{x}_i) - \mathbb{E}_D\left[h\left(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D\right)\right]\right)^2$$

Bias²

$+$

$$+ \sum_i \mathbb{E}_D\left[\left\{h\left(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D\right) - \mathbb{E}_D\left[h\left(\boldsymbol{x}_i; \hat{\boldsymbol{\theta}}_D\right)\right]\right\}^2\right]$$

Variance

sample size dependent!

$+$

$$+ \sum_i \sigma_\epsilon^2$$

Measurement error ( Noise )

# How to quantify expected *out-of-sample* error? Putting them altogether

$$\mathbb{E}_{D,\epsilon}\left[\ell\left(\boldsymbol{y}, h\left(\boldsymbol{X};\hat{\boldsymbol{\theta}}_D\right)\right)\right] = \sum_i \left(f\left(\boldsymbol{x}_i\right) - \mathbb{E}_D\left[h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right]\right)^2 \qquad \text{Bias}^2$$

$$+ \sum_i \mathbb{E}_D\left[\left\{h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right) - \mathbb{E}_D\left[h\left(\boldsymbol{x}_i;\hat{\boldsymbol{\theta}}_D\right)\right]\right\}^2\right] \qquad \text{Variance}$$

sample size dependent!

$$+ \sum_i \sigma_\epsilon^2$$

Measurement error ( Noise )

This picture is valid if the consistency condition holds

$$\mathbb{E}_D\left[h\left(\boldsymbol{x};\hat{\boldsymbol{\theta}}_D\right)\right] \approx \lim_{m\to\infty} h_S^* = h_{\mathcal{D}}^*.$$



$Y^X$  (all functions)

$\mathcal{H}$ (hypotheses space)

$h$

$h_S^*$

$h_{\mathcal{D}}^*$

$y_{\mathcal{D}}^*$

Training Error
(**Optimization**)

Estimation Error
(**Generalization**)
(Variance)

Approximation Error
(**Expressiveness**)
(Bias)

How to quantify expected *out-of-sample* error? Putting them altogether

$$\mathbb{E}_{D,\epsilon}\left[\ell\left(y, h\left(X; \hat{\boldsymbol{\theta}}_D\right)\right)\right] = \sum_i \left(f(x_i) - \mathbb{E}_D\left[h\left(x_i; \hat{\boldsymbol{\theta}}_D\right)\right]\right)^2$$

Bias$^2$

$$+ \sum_i \mathbb{E}_D\left[\left\{h\left(x_i; \hat{\boldsymbol{\theta}}_D\right) - \mathbb{E}_D\left[h\left(x_i; \hat{\boldsymbol{\theta}}_D\right)\right]\right\}^2\right]$$

+

Variance — sample size dependent!

+

$$+ \sum_i \sigma_\epsilon^2$$

Measurement error ( Noise )

$Y^X$ (all functions)

$\mathcal{H}$ (hypotheses space)

$h$

$h_S^*$

$h_{\mathcal{D}}^*$

$y_{\mathcal{D}}^*$

Training Error
(**Optimization**)

Estimation Error
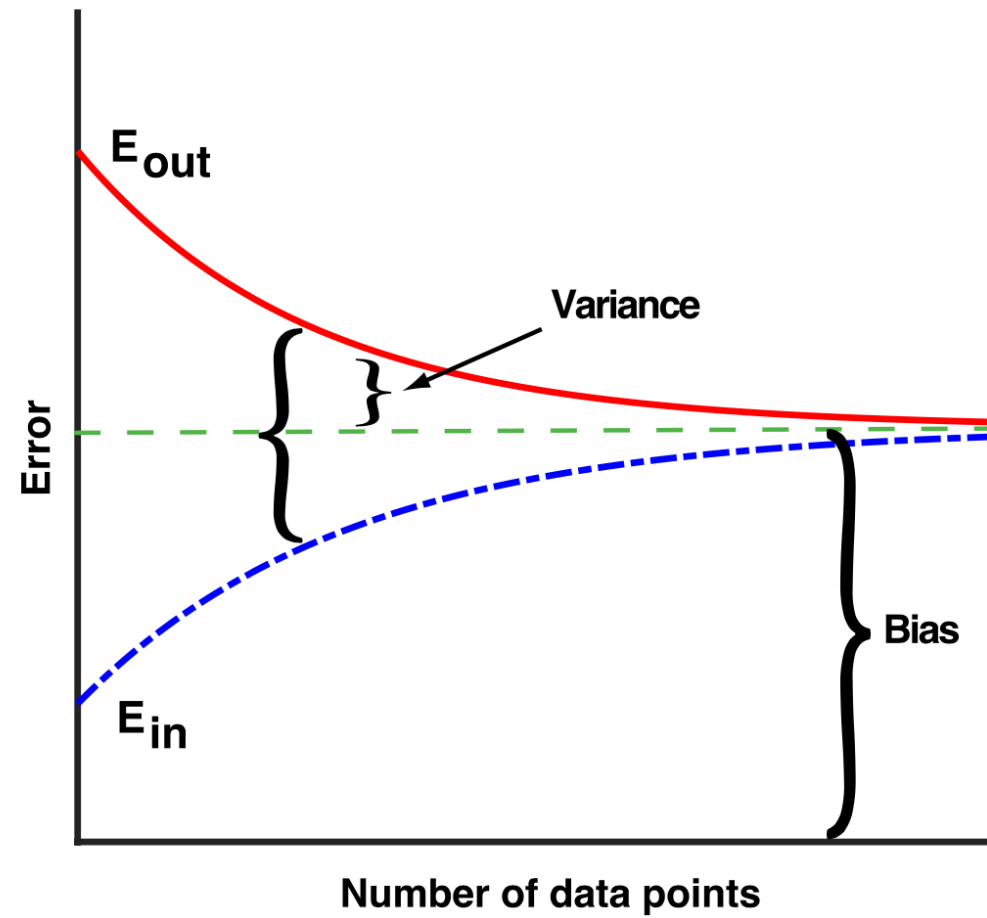(**Generalization**)
(Variance)

Approximation Error
(**Expressiveness**)
(Bias)

This picture is valid if the consistency condition holds

$$\mathbb{E}_D\left[h\left(x; \hat{\boldsymbol{\theta}}_D\right)\right] \approx \lim_{m \to \infty} h_S^* = h_{\mathcal{D}}^*.$$



Error

Optimal Model Complexity

Total Error

Variance

Bias$^2$

Model Complexity

Note: this graph is not quite right…

Measurement error is troublesome!

It makes out-of-sample non-zero even our hypothesis class the same as the god-given signal!

Can we design an alternative approach to learn both signal and noise?

Yes! Bayesian and Probabilistic Inference

Back to a noisy god-given rule

$$y = f(\boldsymbol{x}) + \epsilon \quad \text{with} \quad \epsilon \in \mathcal{N}(0, \sigma_\epsilon^2)$$

We can write this rule for generating the continuous label as

$$P(y \mid \boldsymbol{x}, \boldsymbol{\theta}) = \mathcal{N}\left(f(\boldsymbol{x}), \sigma_\epsilon^2\right)$$

This is a *generative model* (probabilistic model that generates data label, given input).

Back to a noisy god-given rule

$$y = f(\boldsymbol{x}) + \epsilon \quad \text{with} \quad \epsilon \in \mathcal{N}(0, \sigma_\epsilon^2)$$

We can write this rule for generating the continuous label as

$$P(y \mid \boldsymbol{x}, \boldsymbol{\theta}) = \mathcal{N}\left(f(\boldsymbol{x}), \sigma_\epsilon^2\right)$$

This is a *generative model* (probabilistic model that generates data label, given input).

What's the framework to evaluate the hypothesis function?

# Frequentist's Interpretation of Probability

$P(A)$ = *long-run relative frequency* with which A occurs in identical repeats of an experiment. "*A*" restricted to propositions about *random variables*.

Law of large number relates frequency of events in repeated experiments to theoretical probability.

Frequentist's Interpretation of Probability

$P(A)$ = *long-run relative frequency* with which A occurs in identical repeats of an experiment. "*A*" restricted to propositions about *random variables*.

Law of large number relates frequency of events in repeated experiments to theoretical probability.

Bayesian Interpretation of Probability

$P(A|B)$ = a real number measure of *the plausibility of a proposition/hypothesis A*, given (conditional on) *the truth* of the information represented by proposition B. "A" can be *any logical proposition*, not restricted to propositions about random variables.

Both views follow the same mathematical rules of probability theory…

# One useful rule (Bayes' Rule)

$$P(A \mid B) = \frac{P(A, B)}{P(B)} = \frac{P(B \mid A)P(A)}{P(B)}$$

Bayes' rule for evaluating plausibility of a scientific hypothesis, given data.

$$P\left(H_i \mid D, I\right) = \frac{P\left(H_i \mid I\right) P\left(D \mid H_i, I\right)}{P(D \mid I)}$$

$H_i$ = proposition asserting the truth of a hypothesis of interest

$I$ = proposition representing our prior information

$D$ = proposition representing data

$P\left(D \mid H_i, I\right)$ = probability of obtaining data $D$; if $H_i$ and $I$ are true (Likelihood function)

$P\left(H_i \mid I\right)$ = Prior probability of hypothesis

$P\left(H_i \mid D, I\right)$ = Posterior probability of hypothesis

$P(D \mid I) = \sum_i P\left(H_i \mid I\right) P\left(D \mid H_i, I\right)$ is the normalization (tricky to evaluate)

Bayes' rule for evaluating plausibility of a scientific hypothesis, given data.

$$P\left(H_i \mid D, I\right) = \frac{P\left(H_i \mid I\right) P\left(D \mid H_i, I\right)}{P(D \mid I)}$$

$H_i$ = proposition asserting the truth of a hypothesis of interest

$I$ = proposition representing our prior information

$D$ = proposition representing data

$P\left(D \mid H_i, I\right)$ = probability of obtaining data $D$; if $H_i$ and $I$ are true (Likelihood function)

$P\left(H_i \mid I\right)$ = Prior probability of hypothesis

$P\left(H_i \mid D, I\right)$ = Posterior probability of hypothesis

$P(D \mid I) = \sum_i P\left(H_i \mid I\right) P\left(D \mid H_i, I\right)$ is the normalization (tricky to evaluate)

Bayes' rule allows you to evaluate the probability that the hypothesis is true once new data arrives! (How your belief changes depending on the incoming data)