Bias-Variance Tradeoff

Mathematics for Machine Learning

Lecturer: Matthew Wicker

Logistics

- Thank you to all who filled out the survey!
- Our final lecture will be an overview lecture where we connect all of the topics and processes together
- After the final lecture, GTA's will be assigned problem topics to discuss in small groups
- You will get the following materials before next Friday's lecture: Practice Exam, formula sheet, slides containing list of examinable material

Material Covered

Models: Linear models, basis expansion, logistic regression, neural networks, Prob. densities, Bayesian density estimation

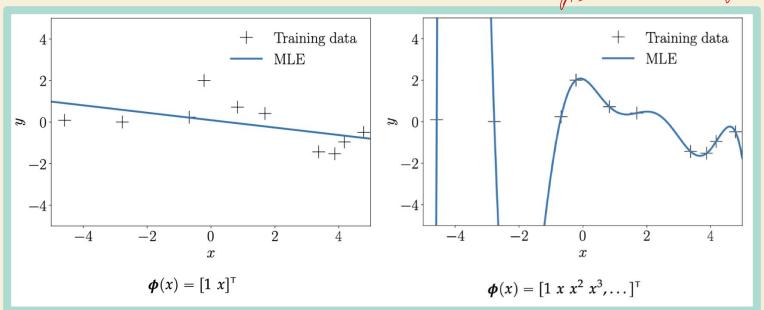
Techniques: Least squares estimation, forward AD, reverse AD, computational graphs, gradient descent, convergence, convexity, Lipschitz continuity, Maximum likelihood, maximum a posteriori, Bayesian inference, LOTUS, change of variables, expectation identities, equating coefficients, joint Gaussian, epistemic/aleatoric uncertainty, concentration inequalities, cross-validation, regularization, generalization error bounds

Settings: Regression, Classification, Density Estimation

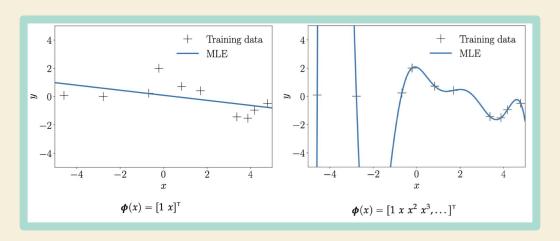
This lecture: Estimators, bias of an estimator, variance of an estimator, bias variance trade off

Recall from last lectures

1. to Selfor Lat worr generalisation



Recall from last lectures



Model validation ideas:

- Test-set -> generalization error w/concentration
- Generalization error bounds for class of models
- Regularization (implicit/explicit)
- Cross-validation

Definition 2.1. Statistic A statistic S is a random variable that is a function of some data \mathcal{D} , $S = g(\mathcal{D})$ where the data \mathcal{D} is a collection of random variables.

An estimator is a statistic that aims at approximating a parameter/property of a distribution.

Definition 2.1. Statistic A statistic S is a random variable that is a function of some data \mathcal{D} , $S = q(\mathcal{D})$ where the data \mathcal{D} is a collection of random variables.

An estimator is a statistic that aims at approximating a parameter/property of a distribution.

sample mean
$$\hat{Z}=rac{Z_1+Z_2+\ldots+Z_N}{N}$$
 of applied to dark

An estimator is a statistic that aims at approximating a parameter/property of a distribution.

eslimb of Z

$$\operatorname{Bias}(\hat{Z}_n) = \mathbb{E}[\hat{Z}_n^{\ell} - Z]$$

Exputed dill between estimper and netral skelistic

An estimator is a statistic that aims at approximating a parameter/property of a distribution.

$$\operatorname{Bias}(\hat{Z}_n) = \mathbb{E}[\hat{Z}_n - Z] = 0$$

Unbiased estimator

We proved unbiased estimator for the sample mean

From lecture 10:

$$\hat{Z} = \frac{Z_1 + Z_2 + \ldots + Z_N}{N}$$

$$\mathbb{E}[\hat{Z}] = rac{\mathbb{E}[Z_1] + \mathbb{E}[Z_2] + \ldots + \mathbb{E}[Z_N]}{N}$$

$$= rac{\mu_1 + \mu_2 + \ldots \mu_N}{N} = \mu$$

$$\hat{\sigma}^2_n = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 \qquad \hat{\sigma}^2_n = \frac{1}{n-1} \sum_{i=1}^n (X_i - \mu)^2$$

Here we have two sample variance estimators, but which one is the biased estimator and which is unbiased?

Letting:
$$\mu=\hat{\mu}_n$$

Let's check the simpler looking equation for bias

$$\hat{\sigma}^2_n = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$$

$$\operatorname{bias}(\hat{\sigma^2}_n) = \mathbb{E}[\hat{\sigma^2}_n - \sigma^2]$$

$$\operatorname{Rim}: \left\{ \left[\left(\widehat{S} - \mathcal{S} \right) \right] \right\} \qquad \operatorname{for}: \left\{ \left[\left(\widehat{S} - \mathcal{F} \left(\mathcal{S} \right) \right)^{2} \right] \right\} \qquad = \operatorname{cor} \left\{ \operatorname{sin}_{\mathcal{S}^{2}} \right\}$$

$$= \left[\left(\widehat{S} \right) \cdot \mathcal{E} \left(\mathcal{S} \right) : \mathcal{E} \left(\widehat{S} \right) \cdot \mathcal{S} \right] \qquad = \left[\left(\left(\widehat{S} - \mathcal{F} \left(\mathcal{S} \right) \right)^{2} \right) \left[\left(\mathcal{S} - \mathcal{S} \right)^{2} \right] \right\} \qquad = \left[\left(\left(\widehat{S} - \mathcal{F} \left(\mathcal{S} \right) \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \right] \qquad = \left[\left(\left(\widehat{S} - \mathcal{F} \left(\mathcal{S} \right) \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \right] \qquad = \left[\left(\left(\widehat{S} - \mathcal{F} \left(\mathcal{S} \right) \right)^{2} \right] \qquad = \left[\left(\left(\widehat{S} - \mathcal{F} \left(\mathcal{S} \right) \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \right] \qquad = \left[\left(\left(\widehat{S} - \mathcal{F} \left(\mathcal{S} \right) \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \right] \qquad = \left[\left(\left(\widehat{S} - \mathcal{F} \left(\mathcal{S} \right) \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \right] \qquad = \left[\left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \right] \qquad = \left[\left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \right] \qquad = \left[\left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \right] \qquad = \left[\left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \right] \qquad = \left[\left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \right] \qquad = \left[\left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \right] \qquad = \left[\left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \right] \qquad = \left[\left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \right] \qquad = \left[\left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \right] \qquad = \left[\left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \right] \qquad = \left[\left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \right] \qquad = \left[\left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \right] \qquad = \left[\left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \left(\left(\mathcal{S} - \mathcal{S} \right)^{2} \right) \right] \qquad = \left[\left($$

Let's check the simpler looking equation for bias

$$\hat{\sigma}^2_n = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$$

$$\text{bias}(\hat{\sigma}^2_n) = \mathbb{E}[\hat{\sigma}^2_n - \sigma^2]$$

$$= \mathbb{E}[\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - \sigma^2]$$

Let's check the simpler looking equation for bias

$$\begin{aligned} \operatorname{bias}(\hat{\sigma^2}_n) &= \mathbb{E}[\hat{\sigma^2}_n - \sigma^2] \\ &= \mathbb{E}[\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - \sigma^2] \quad \text{for } \sigma^2 \\ &= \mathbb{E}[\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2] - \sigma^2 \end{aligned}$$

Simplifying the expression in the expectation:

$$\frac{1}{n}\sum_{i=1}^{n}(X_i-\mu)^2 = \frac{1}{n}\sum_{i=1}^{n}(X_i^2-2X_i\mu+\mu^2)$$

Simplifying the expression in the expectation:

$$\frac{1}{n} \sum_{i=1}^{n} (X_i - \mu)^2 = \frac{1}{n} \sum_{i=1}^{n} (X_i^2 - 2X_i \mu + \mu^2)$$

$$\propto \sum_{i=1}^{n} X_i^2 - \sum_{i=1}^{n} 2X_i \mu + \sum_{i=1}^{n} \mu^2$$

Simplifying the expression in the expectation:

$$\frac{1}{n} \sum_{i=1}^{n} (X_i - \mu)^2 = \frac{1}{n} \sum_{i=1}^{n} (X_i^2 - 2X_i \mu + \mu^2)$$

$$\propto \sum_{i=1}^{n} X_i^2 - \sum_{i=1}^{n} 2X_i \mu + \sum_{i=1}^{n} \mu^2$$

$$= \left(\sum_{i=1}^{n} X_i^2\right) - 2\mu \left(\sum_{i=1}^{n} X_i\right) + n\mu^2$$

Simplifying the expression in the expectation:

$$= \left(\sum_{i=1}^{n} X_i^2\right) - 2\mu \left(\sum_{i=1}^{n} X_i\right) + n\mu^2$$

$$= \left(\sum_{i=1}^{n} X_i^2\right) - 2n\mu^2 + n\mu^2$$

$$= \left(\sum_{i=1}^{n} X_i^2\right) - 2n\mu^2 + n\mu^2$$

$$= \left(\sum_{i=1}^{n} X_i^2\right) - 2n\mu^2 + n\mu^2$$

Simplifying the expression in the expectation:

$$= \left(\sum_{i=1}^{n} X_i^2\right) - 2\mu \left(\sum_{i=1}^{n} X_i\right) + n\mu^2$$

$$= \left(\sum_{i=1}^{n} X_i^2\right) - 2n\mu^2 + n\mu^2$$

$$= \left(\sum_{i=1}^{n} X_i^2\right) - n\mu^2$$

Simplifying the expression in the expectation:

$$\mathbb{E}\left[\left(\sum_{i=1}^{n} X_i^2\right) - n\mu^2\right] = \left(\sum_{i=1}^{n} \mathbb{E}[X_i^2]\right) - \mathbb{E}[n\mu^2]$$

Exactadion over luka (xi)

Simplifying the expression in the expectation:

$$\mathbb{E}\left|\left(\sum_{i=1}^{n} X_i^2\right) - n\mu^2\right| = \left(\sum_{i=1}^{n} \mathbb{E}[X_i^2]\right) - \mathbb{E}[n\mu^2]$$

Aside, identities we will prove later:

Hint: this is just rearranging the variance definition

$$\mathbb{E}[X_i^2] = \sigma^2 + \mu^2$$

$$\mathbb{E}[\mu^2] = \frac{\sigma^2}{n} + \mu^2$$

Simplifying the expression in the expectation:

$$= \left(\sum_{i=1}^{n} \mathbb{E}[X_i^2]\right) - \mathbb{E}[n\mu^2]$$

$$= \sum_{i=1}^{n} (\sigma^2 + \mu^2) - n \left(\frac{\sigma^2}{n} + \mu^2 \right)$$

$$= n\sigma^2 + n\mu^2 - \sigma^2 - n\mu^2$$

$$=(n-1)\sigma^2$$

Simplifying the expression in the expectation:

$$= \left(\sum_{i=1}^{n} \mathbb{E}[X_i^2]\right) - \mathbb{E}[n\mu^2]$$

$$= n\sigma^2 + n\mu^2 - \sigma^2 - n\mu^2$$

$$= (n-1)\sigma^2$$

$$\operatorname{bias}(\hat{\sigma^2}_n) = \frac{n-1}{n}\sigma^2 - \sigma^2$$

$$\hat{\sigma}^2_n = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$$

By plugging in our estimator into the definition of bias, making the i.i.d. assumption and being comfortable with manipulating expectation identities, we showed the bias is non-zero:

$$\operatorname{bias}(\hat{\sigma^2}_n) = \frac{n-1}{n}\sigma^2 - \sigma^2$$

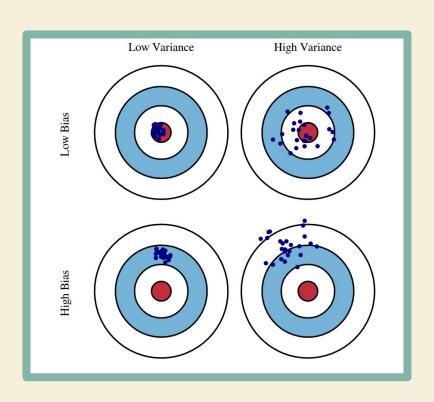
Variance of an estimator

Definition 2.1. Statistic A statistic S is a random variable that is a function of some data \mathcal{D} , $S = g(\mathcal{D})$ where the data \mathcal{D} is a collection of random variables.

An estimator is a statistic that aims at approximating a parameter/property of a distribution.

$$\operatorname{Var}(\hat{S}) = \mathbb{E}\left[(\hat{S} - \mathbb{E}(\hat{S}))^2\right]$$

Visualization of bias-variance



Bias-variance decomposition

$$\operatorname{Err}(\hat{S}) = \mathbb{E}[(\hat{S} - S)^2] = \operatorname{Bias}^2 + \operatorname{Var}_{\text{eccord}}$$

The bias-variance decomposition simply tells us that the (mean squared error) of an estimator is completely accounted for by the sum of the squared bias and variance of the estimator.

Bias-variance decomposition

$$\operatorname{Err}(\hat{S}) = \mathbb{E}[(\hat{S} - S)^2] = \operatorname{Bias}^2 + \operatorname{Var}$$

$$= \mathbb{E}[(g(\mathcal{D}) - S)^2] \qquad \text{only when for most of with the following of the property of the propert$$

The bias-variance decomposition simply tells us that the (mean squared error) of an estimator is completely accounted for by the sum of the squared bias and variance of the estimator.

6'05: EL S-5] UN: EL(ECS)-5/3

$$\mathbb{E}[(\hat{S} - S)^2] = \mathbb{E}\left[\left(\hat{S} - \mu + \mu - S\right)^2\right] \qquad \stackrel{\text{t. }}{=} \mathbb{E}\left[\left((\hat{S} - \mu) + (\mu - S)\right)^2\right]$$

-
$$\mathcal{E}\left(\hat{S}-\mathcal{E}(\hat{S})^2\right)$$
 + $\mathcal{E}\left(\mathcal{E}(\hat{S})-\hat{S}^2\right)$
Letting: $\mu=\mathbb{E}[\hat{S}]$

$$\mathbb{E}[(\hat{S} - S)^2] = \mathbb{E}\left[\left(\hat{S} - \mu + \mu - S\right)^2\right]$$

$$= \mathbb{E}\left[\left((\hat{S} - \mu) + (\mu - S)\right)^2\right]$$

$$= \mathbb{E}\left[\left((\hat{S} - \mu)^2 + 2(\hat{S} - \mu)(\mu - S) + (\mu - S)^2\right]\right]$$

$$\begin{split} \mathbb{E}[(\hat{S}-S)^2] &= \mathbb{E}\left[\left(\hat{S}-\mu+\mu-S\right)^2\right] \\ &= \mathbb{E}\left[\left((\hat{S}-\mu)+(\mu-S)\right)^2\right] \\ &= \mathbb{E}\left[\left(\hat{S}-\mu\right)^2+2(\hat{S}-\mu)(\mu-S)+(\mu-S)^2\right] \\ &= \mathbb{E}\left[(\hat{S}-\mu)^2+(\mu-S)^2\right] \quad (\star) \end{split}$$

$$\mathbb{E}[(\hat{S} - S)^2] = \mathbb{E}\left[\left(\hat{S} - \mu + \mu - S\right)^2\right]$$

$$= \mathbb{E}\left[\left((\hat{S} - \mu) + (\mu - S)\right)^2\right]$$

$$= \mathbb{E}\left[\left((\hat{S} - \mu) + (\mu - S)\right)^2\right]$$

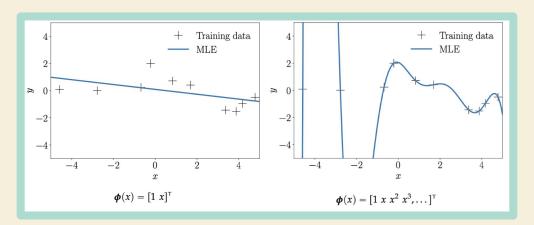
$$= \mathbb{E}\left[(\hat{S} - \mu)^2 + 2(\hat{S} - \mu)(\mu - S) + (\mu - S)^2\right]$$

$$= \mathbb{E}\left[(\hat{S} - \mu)^2 + (\mu - S)^2\right] \quad (\star)$$

$$= \mathbb{E}\left[(\hat{S} - \mu)^2\right] + \mathbb{E}\left[(\mu - S)^2\right]$$

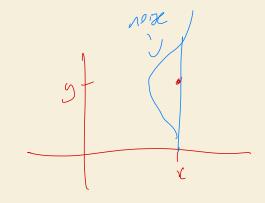
$$= \mathbb{E}\left[(\hat{S} - \mu)^2\right] + (\mu - S)^2 \quad \text{for } s = 1 \text{ for } s = 1 \text{ fo$$

Our course notes discuss the particulars of bias-variance decomposition in linear models, which you should study (in order to answer exercises), but in lecture we will be more general/conceptual.



https://www.cs.cornell.edu/courses/cs4780/2022sp/notes/LectureNotes17.html

$$\hat{y} = \int_{y} P(y|\mathbf{x}) dy$$



Consider the true label of our machine learning models to be the integral over the likelihood

$$\hat{y} = \int_y P(y|\mathbf{x}) dy$$

We can write the error of our model as the integral over both the output and input space of our mean squared error:

$$\mathbb{E}_{P(\mathbf{x},y)}[(f^{\theta}(\mathbf{x}) - y)^{2}] = \int_{x} \int_{y} (f^{\theta}(\mathbf{x}) - y)^{2} P(\mathbf{x},y) dy d\mathbf{x}$$

$$\mathbb{E}_{P(\mathbf{x},y)}[(f^{\theta}(\mathbf{x}) - y)^{2}] = \int_{x} \int_{y} (f^{\theta}(\mathbf{x}) - y)^{2} P(\mathbf{x},y) dy d\mathbf{x}$$

When we compute the bias and variance of an estimator (in this case the estimator of our risk of a model), we need to do so with respect to the "average model"

$$\hat{f}^{ heta} = \mathbb{E}_{\mathcal{D} \sim P(\mathbf{x}, y)} \int_{\mathcal{D}} f^{ heta}_{\mathcal{D}} P(\mathcal{D}) d\mathcal{D}$$

are nowld sixty integrate on all potential datasets

Now, we have defined all of the components we need to quantify the bias and variance of our model error estimator:

$$\mathbb{E}_{\mathbf{x},y,\mathcal{D}}\left[\left(f_{\mathcal{D}}^{\theta}(\mathbf{x})-y\right)^{2}\right] \sim \left\{\left(\left(\hat{\mathbf{x}}-\mathbf{E}(\mathbf{x})\right)^{2}\right) + \mathbf{E}\left(\mathbf{E}(\mathbf{x})-\mathbf{x}\right)^{2}\right\}$$

$$= \left(\left(\left(f_{\mathcal{D}}^{\theta}(\mathbf{x})-\hat{\mathbf{y}}\right)^{2}\right) + \mathbf{E}\left(\left(f_{\mathcal{D}}^{\theta}(\mathbf{x})-\mathbf{y}\right)^{2}\right)^{2}$$

$$\mathbb{E}_{\mathbf{x},y,\mathcal{D}}\left[(f_{\mathcal{D}}^{\theta}(\mathbf{x})-y)^{2}\right]$$

As before, we consider subtracting and adding the mean from the estimated quantity:

$$=\mathbb{E}_{\mathbf{x},y,\mathcal{D}}\left[\left((f_{\mathcal{D}}^{ heta}(\mathbf{x})-\hat{f}^{ heta}(\mathbf{x}))
ight. + \left(\hat{f}^{ heta}(\mathbf{x})-y
ight)^{2}
ight]$$

By the same cancellation logic, we get that the middle terms of completing the square drop off, so:

$$= \mathbb{E}_{\mathbf{x},y,\mathcal{D}} \left[\left((f_{\mathcal{D}}^{\theta}(\mathbf{x}) - \hat{f}^{\theta}(\mathbf{x}))^2 + (\hat{f}^{\theta}(\mathbf{x}) - y)^2 \right) \right]$$

And then by linearity of expectation we have the definition of bias-variance again:

$$\mathbb{E}_{\mathbf{x},\mathcal{D}}\left[(f_{\mathcal{D}}^{\theta}(\mathbf{x}) - \hat{f}^{\theta}(\mathbf{x}))^{2}\right] + \mathbb{E}_{\mathbf{x},y}[(\hat{f}^{\theta}(\mathbf{x}) - y)^{2}]$$

for model derivation and with definition and the first definition of the property of the property

$$\mathbb{E}_{\mathbf{x},\mathcal{D}}\left[(f_{\mathcal{D}}^{\theta}(\mathbf{x}) - \hat{f}^{\theta}(\mathbf{x}))^{2}\right] + \mathbb{E}_{\mathbf{x},y}\left[(\hat{f}^{\theta}(\mathbf{x}) - y)^{2}\right]$$

Notice! Another MSE term here, this time the risk of our average model. We can repeat our bias-variance steps to decompose this further

Again, we can add and subtract the estimator for y to further decompose this term:

$$\mathbb{E}_{\mathbf{x},\mathcal{D}}\left[(f_{\mathcal{D}}^{\theta}(\mathbf{x}) - \hat{f}^{\theta}(\mathbf{x}))^{2}\right] + \mathbb{E}_{\mathbf{x},y}[(\hat{f}^{\theta}(\mathbf{x}) - y)^{2}]$$

$$\mathbb{E}_{\mathbf{x},y}[(\hat{f}^{\theta}(\mathbf{x}) - y)^2] = \mathbb{E}_{\mathbf{x},y}\left[(\hat{f}^{\theta}(\mathbf{x}) - \hat{y} + \hat{y} - y)^2\right]$$

Again, we can add and subtract the estimator for y to further decompose this term:

$$\mathbb{E}_{\mathbf{x},\mathcal{D}}\left[(f_{\mathcal{D}}^{\theta}(\mathbf{x}) - \hat{f}^{\theta}(\mathbf{x}))^{2}\right] + \mathbb{E}_{\mathbf{x},y}\left[(\hat{f}^{\theta}(\mathbf{x}) - y)^{2}\right]$$

$$\mathbb{E}_{\mathbf{x},y}[(\hat{f}^{\theta}(\mathbf{x}) - y)^{2}] = \mathbb{E}_{\mathbf{x},y}\left[(\hat{f}^{\theta}(\mathbf{x}) - \hat{y} + \hat{y} - y)^{2}\right]$$

$$= \mathbb{E}_{\mathbf{x},y}[(\hat{y} - y)^{2}] + \mathbb{E}_{\mathbf{x}}\left[(\hat{f}^{\theta}(\mathbf{x}) - \hat{y})^{2}\right]$$

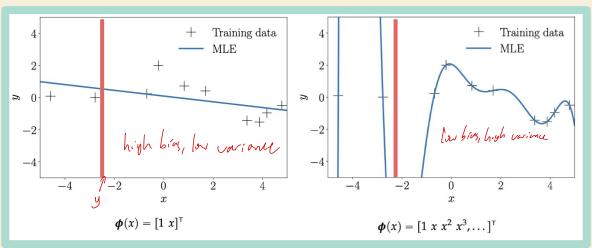
Final bias-variance decomposition

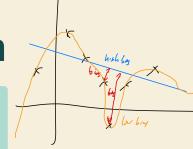
$$\mathbb{E}_{\mathbf{x},y,\mathcal{D}}\left[(f_{\mathcal{D}}^{\theta}(\mathbf{x}) - y)^2 \right]$$

We decomposed the error into three terms:

$$\mathbb{E}_{\mathbf{x},y}[(\hat{y}-y)^2] + \mathbb{E}_{\mathbf{x}} \quad [(\hat{f}^{\theta}(\mathbf{x})-\hat{y})^2] + \mathbb{E}_{\mathbf{x},\mathcal{D}}[(\hat{f}^{\theta}_{\mathcal{D}}(\mathbf{x})-\hat{f}^{\theta}(\mathbf{x}))^2]$$
Noise Squared bias Variance

Visualizing this decomposition





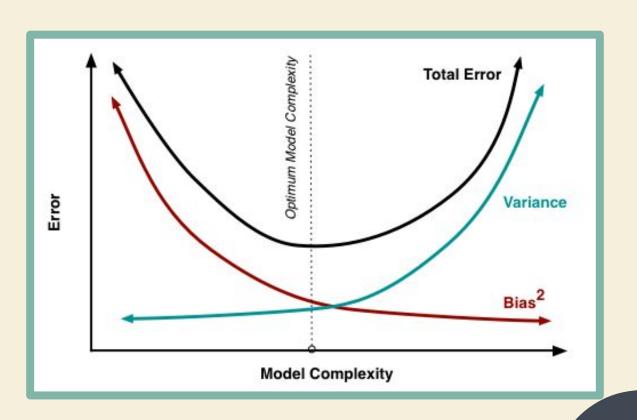
$$\mathbb{E}_{\mathbf{x},y}[(\hat{y}-y)^2] + \mathbb{E}_{\mathbf{x}} [(\hat{f}^{\theta}(\mathbf{x}) - \hat{y})^2] + \mathbb{E}_{\mathbf{x},\mathcal{D}}[(\hat{f}^{\theta}_{\mathcal{D}}(\mathbf{x}) - \hat{f}^{\theta}(\mathbf{x}))^2]$$

Noise

Squared bias

Variance

Model selection



Next lecture: Principal Component-Analysis

