

INF8953 CE - Fall 2020

Machine Learning

- Sarath Chandar

4. Empirical Risk Minimization.



4. Empirical Risk Minimization

Let $x \in \mathbb{R}^p$, $y \in \mathbb{R}$

$h: x \rightarrow y$ [hypothesis]

When we use model based algorithms, model defines a function 'f' parameterized by 'w'. For every possible parameter configuration w^+ , $f(w^+)$ is a hypothesis. Finding right set of parameters for the model is equivalent to finding the right hypothesis from the given set of hypothesis supported by f.

$$EPE(h) = E[L(y, h(x))]$$

$$= \iint L(y, h(x)) p(x, y) dx dy$$

This is also known as the risk associated with the hypothesis $h: R(h)$

Goal: Find a hypothesis h^* among the class of functions for which the risk $R(h)$ is minimal.

$$h^* = \underset{h \in \mathcal{H}}{\operatorname{argmin}} R(h)$$

In general, $R(h)$ cannot be computed because the distribution $P(x, y)$ is unknown to the learning algorithm.

However, we can compute an approximation for this risk by using the training set.

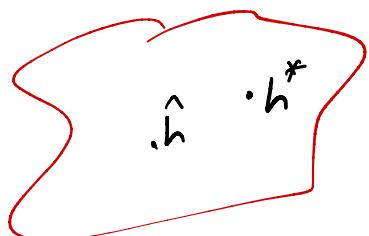
$$R_{\text{emp}}(h) = \frac{1}{N} \sum_{i=1}^N L(y^{(i)}, h(x^{(i)}))$$

empirical risk.

ERM: find \hat{h} which minimizes the empirical risk.

$$\hat{h} = \underset{h \in \mathcal{H}}{\operatorname{argmin}} R_{\text{emp}}(h)$$

Note 1: There is no guarantee that $h^* = \hat{h}$.



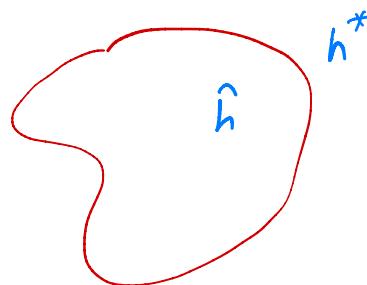
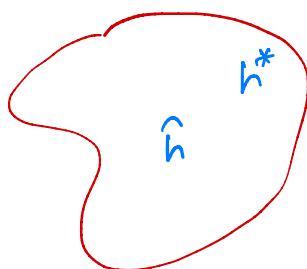
Note 2: The inductive bias of the algorithm restricts the hypothesis class.

ex: linear models restrict the hypothesis class to linear functions only.

Original hypothesis class: \mathcal{H} subset
restricted hypothesis class: $\mathcal{H}' \subseteq \mathcal{H}$

The restricted hypothesis class may or may not contain h^* .

Case 1: $h^* \in \mathcal{H}'$ Case 2: $h^* \notin \mathcal{H}'$



Bias - Variance decomposition :-

Consider regression problem:

$$L(y, h(x)) = (h(x) - y)^2$$

$$h^* = \underset{h \in \mathcal{H}}{\operatorname{argmin}} E[L(y, h(x))]$$

$$\boxed{h^* = E[y|x]} \quad \text{regression function.}$$

$$\begin{aligned} \{h(n) - y\}^2 &= \left\{ h(n) - \underbrace{E[y|x]}_{\text{regression function}} + \underbrace{E[y|x] - y}_{\text{error term}} \right\}^2 \\ &= (h(n) - E[y|x])^2 + (E[y|x] - y)^2 \\ &\quad + 2(h(n) - E[y|x])(E[y|x] - y) \end{aligned}$$

$$\begin{aligned} E[L(y, h(x))] &= E[\{h(n) - y\}^2] \\ &= E[(h(n) - E[y|x])^2] + E[(E[y|x] - y)^2] \\ &\quad + 2 E[(h(n) - E[y|x])(E[y|x] - y)] \\ &\quad \text{③} \end{aligned}$$

First, we will show that ③ = 0.

$$\begin{aligned} ③ &= E_x E_{y|x} (h(n) E[y|x] - h(n) y - E[y|x]^2 + y E[y|x]) \\ &= E_x (h(x) E[y|x] - h(n) E[y|x] - E[y|x]^2 + E[y|x]^2) \\ &= E_x [0] = 0. \end{aligned}$$

$$E[L(y, h(x))] = \underbrace{E[(h(x) - E[y|x])^2]}_{\textcircled{1}} + \underbrace{E[(E[y|x] - y)^2]}_{\textcircled{2}}$$

$$\textcircled{1} E[(h(x) - E[y|x])^2]$$

$$= \int_x \int_n \underbrace{(h(x) - E[y|x])^2}_{\text{Not a function of } y} p(y|x) dy p(n) dn$$

$$= \int_n (h(x) - E[y|x])^2 p(n) dn$$

$$\textcircled{2} E[(E[y|x] - y)^2]$$

$$= \int_n \int_{y|x} \underbrace{(y - E[y|x])^2}_{\text{Var}(y|x)} p(y|x) dy p(n) dn$$

$$= \int_n \text{Var}(y|x) p(n) dn.$$

$$E[L(y, h(x))] = E[\{h(x) - y\}^2]$$

$$= E_x[(h(x) - E[y|x])^2]$$

$$+ E_x[\text{Var}(y|x)]$$

Note 1: Second term is the variance of the distribution y , averaged over n . It represents the intrinsic variable of the target data and can be regarded as noise.

Second term is independent of $h(n)$

if $h(n) = E[y|x]$,
then $\textcircled{1} \Rightarrow$

↳ irreducible minimum value of the loss fn.

Note 2: First term - is a function of $h(n)$.

This is minimum only when $h(n) = E[y|x]$.

This 0 happens only when we know the true function. Even Equ(1) is zero then also we have some irreducible error(variance-noise)

Consider : $(h(n) - E[y|x])^2$

In practice, we have a dataset \mathcal{D} containing only a finite number N of data points.

For any dataset \mathcal{D} , we run our algorithm and find $h(n)$. Call it $h(n; \mathcal{D})$.

$$\{ h(n; \mathcal{D}) - E[y|x] \}^2$$

$$\begin{aligned}
 & E_D \left[\{h(n; D) - E[y|x]\}^2 \right] \\
 &= E_D \left[\underbrace{\{h(n; D) - E_D[h(n; D)]\}}_{+} + \underbrace{E_D[h(n; D)] - E[y|x]}_{-} \right]^2 \\
 &= E_D \left[\{h(n; D) - E_D[h(n; D)]\}^2 + \{E_D[h(n; D)] - E[y|x]\}^2 \right. \\
 &\quad \left. + 2(h(n; D) - E_D[h(n; D)])(E_D[h(n; D)] - E[y|x]) \right] \\
 &= 0 \quad [\text{Exercise}] \\
 &= \left[E_D[h(n; D)] - E[y|x] \right]^2 + E_D \left[\{h(n; D) - E_D[h(n; D)]\}^2 \right] \\
 &= \text{bias}^2 + \text{variance}
 \end{aligned}$$

bias - extent to which the average prediction
 over all datasets vary around the average.

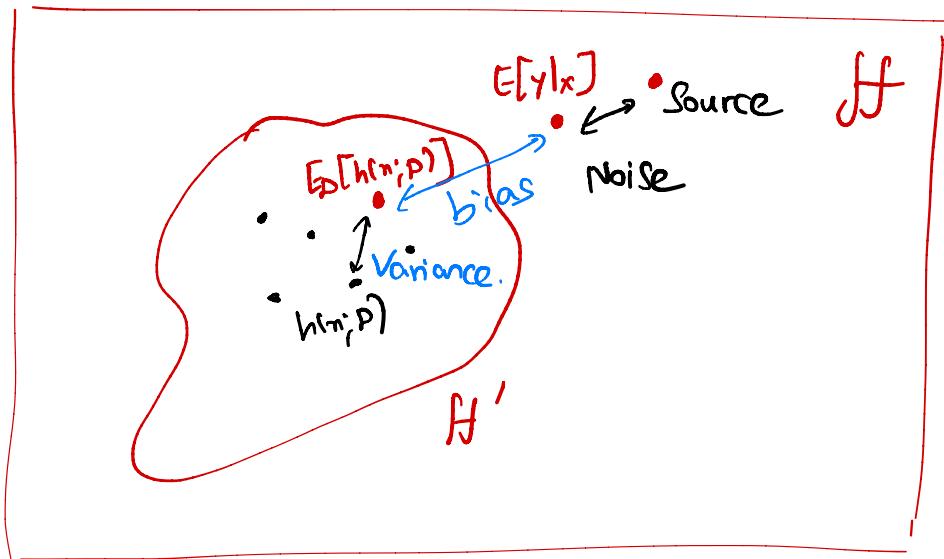
Variance - extent to which solutions for individual datasets vary around their average,
 ↳ extent to which $h(n; D)$ is sensitive to particular choice of dataset.

Expected loss = $(\text{bias})^2 + \text{Variance} + \text{noise}$.

$$(\text{bias})^2 = \int \left\{ E_{\mathcal{D}}[h(x; p)] - E[y|x] \right\}^2 P_{xy} dx$$

$$\text{Variance} = \int \left\{ E_{\mathcal{D}}[h(x; p)] - E_{\mathcal{D}}[h(x; p)] \right\}^2 P_{yy} dy$$

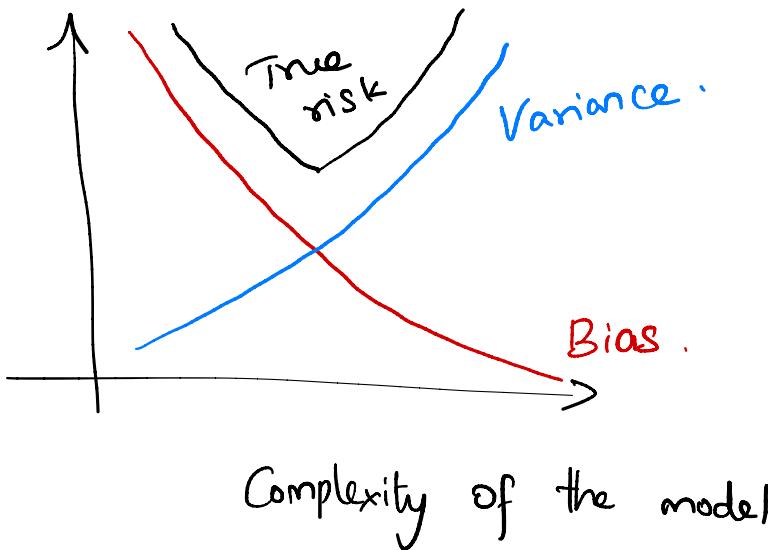
$$\text{Noise} = \int \left\{ E[y|x] - y \right\}^2 P_{(x,y)} dx dy$$



Note: There is a tradeoff between bias and variance.

Very flexible models — low bias and high variance.

Rigid models — high bias and low variance.



Occam's Razor:-

"One should not increase, beyond what is necessary, the number of entities required to explain anything".

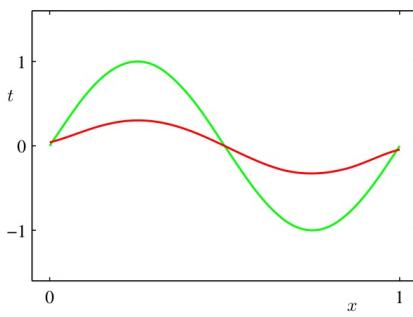
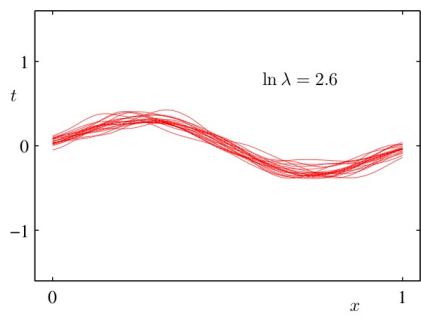
"Seek the simplest explanation".

Regression example:-

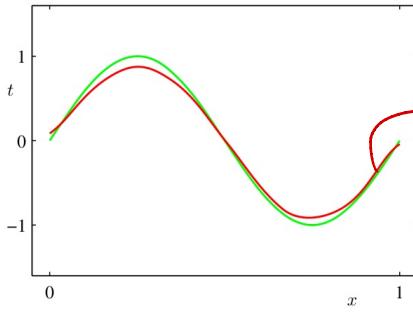
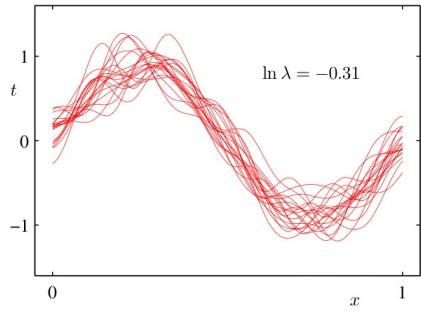
100 datasets each containing $N = 25$ data points from $\sin(2\pi n)$ curve.

$$L = \# \text{ of datasets} = 100$$

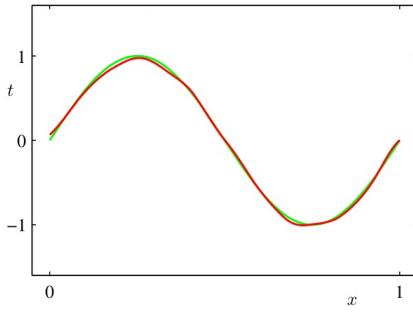
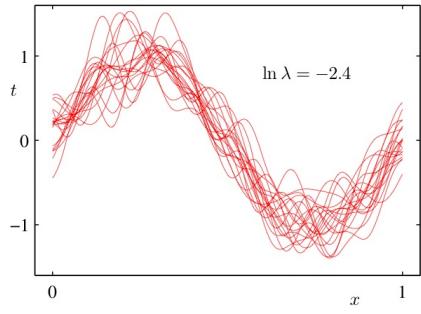
Model: 24 Gaussian basis functions.



λ is large.
variance \rightarrow low
bias \rightarrow high



red line
↳ average fit.

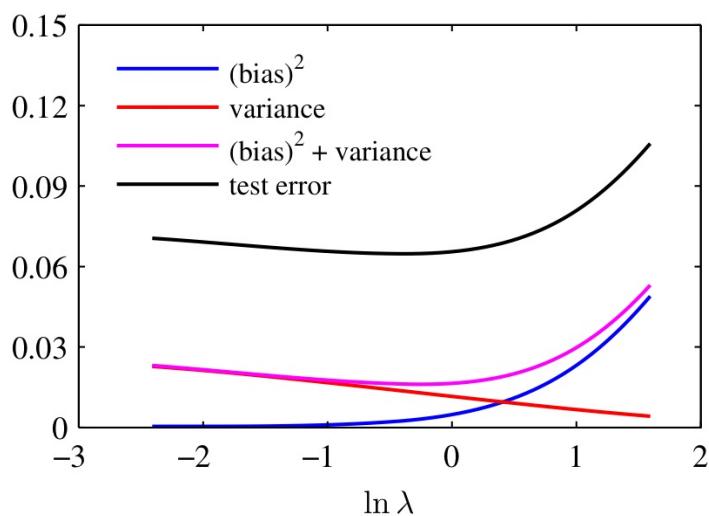


λ is small.
variance \rightarrow high
bias \rightarrow low.

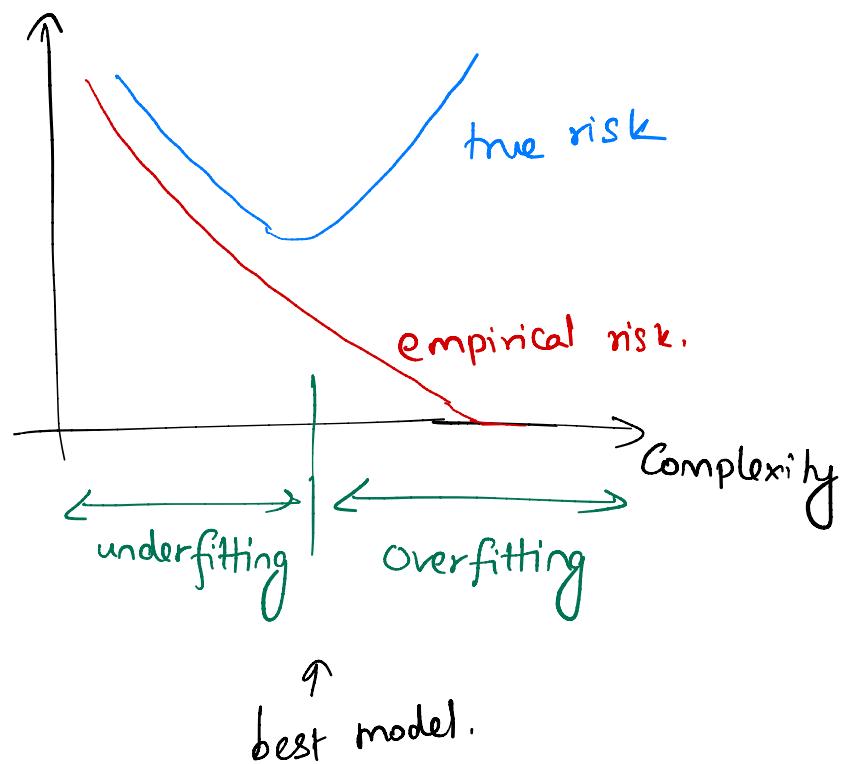
$$\text{Average prediction : } \bar{h}(n) = \frac{1}{L} \sum_{l=1}^L h^{(l)}(x)$$

$$(\text{bias})^2 = \frac{1}{N} \sum_{n=1}^N \left\{ \bar{h}(n^{(n)}) - E[y^{(n)}|x^{(n)}] \right\}^2$$

$$\text{Variance} = \frac{1}{N} \sum_{n=1}^N \frac{1}{L} \sum_{l=1}^L \left\{ h^{(l)}(x^{(n)}) - \bar{h}(n^{(n)}) \right\}^2$$



Relation between true risk, empirical risk, overfitting :-



You should know!

1. Hypothesis
2. True risk
3. Empirical risk.
4. Bias-variance decomposition.
5. Irreducible error.
6. Occam's razor.

]. underfitting / overfitting / true risk / empirical
risk / Complexity of the model - Connections .
