# Lecture 8

B & V for Kernel Regression

Bias Case $\overline{II}$

Overfitting, Underfitting, B-V tradeoff

Neural Networks (start)

# Lecture II: Prediction – Basic concepts

Marina Meilă
mmp@stat.washington.edu

Department of Statistics
University of Washington

October, 2023

Parametric vs non-parametric

**Generative and discriminative models for classification** ✔
   Generative classifiers
   Discriminative classifiers
   Generative vs discriminative classifiers

Loss functions ╱
   Bayes loss

Variance, bias and complexity
   Case 1: With $f^{true}$, Least Squares loss
   Bias and Variance for Kernel Regression (Lecture II.1) ⬅
   Case 2: Bias as as model (mis)fit ⬅

Sampling Variance

**Reading** HTF Ch.: 2.1–5,2.9, 7.1–4 bias-variance tradeoff, Murphy Ch.: 1., 8.6[1], Bach Ch.:

─────────────────────────

[1]Neither textbook is close to these notes except in a few places; take them as alternative perspectives or related reading

# The "learning" problem

▶ **Given**
▶ a problem (e.g. recognize digits from $m \times m$ gray-scale images)
▶ a **sample** or (**training set**) of **labeled data**

$$\mathcal{D} = \{(x^1, y^1), (x^2, y^2), \dots (x^n, y^n)\}$$

drawn i.i.d. from an unknown $P_{XY}$
▶ **model class** $\mathcal{F} = \{f\}$ = set of predictors to choose from

▶ **Wanted**
▶ a predictor $f \in \mathcal{F}$ that performs well on future samples from the same $P_{XY}$

  ▶ "choose a predictor $f \in \mathcal{F}$" = training/learning
  ▶ "performs well on future samples" (i.e. $f$ **generalizes** well) – how do we measure this? how can we "guarantee" it?
  ▶ choosing $\mathcal{F}$ is the **model selection problem** – about this later

# Bias and variance: Preliminaries

## Setup/What we have

- ▶ a data source $P_{XY}$
- ▶ a class of predictors $\mathcal{F}$
- ▶ From $P_{XY}$ we sample i.i.d. $\mathcal{D}_n$ of size $n$. Hence $\mathcal{D}_n \sim P_{XY}^n$.
- ▶ A training algorithm that estimates/chooses/learns $\hat{f}_n$ from $\mathcal{D}_n$.
    - ▶ minimize $\hat{L}_{f \in \mathcal{F}}$ (empirical loss) Example CART, Logistic Regression and all Max Likelihood methods
    - ▶ minimize over $f \in \mathcal{F}$ (**regularized** loss)

$$\hat{L}(f) + \lambda J(f) \tag{15}$$

    with $\lambda > 0$ a regularization parameter Example Ridge Regression, SVM
    - ▶ other training method (e.g. K-NN, LDA)

# Bias and Variance for parameter estimation

▶ We want to estimate a parameter $\theta \in \Theta \subseteq \mathbb{R}$

▶ We use $\mathcal{D}_n$ to obtain estimator $\hat{\theta}_{\mathcal{D}_n}$ which is a function of $\mathcal{D}_n$.

▶ $\mathcal{D}_n$ is random, hence so is $\hat{\theta}_{\mathcal{D}_n}$.

▶ Bias$= (\hat{\theta}_{\mathcal{D}_n}) = E_{P^n}[\hat{\theta}_{\mathcal{D}_n}] - \theta$

▶ Variance$= Var_{P^n}(\hat{\theta}_{\mathcal{D}_n})$
   Both Bias and Variance are computed under the distribution from which we sampled $\mathcal{D}_n$, denoted by $P^n$.

▶ Example Estimating $\mu, \sigma^2$ for $N(\mu, \sigma^2)$, $\mathcal{D}_n = \{x_{1:n}\} \subset \mathbb{R}$

$$\hat{\mu} = \frac{1}{n}\sum_i x^i \tag{16}$$

$$Bias(\hat{\mu}) = E[\hat{\mu} - \mu] = \mu - \mu = 0 \quad \hat{\mu} \text{ is unbiased} \tag{17}$$

$$Var(\hat{\mu}) = \sigma^2/n \tag{18}$$

$$\hat{\sigma}^2 = \frac{1}{n}\sum_i (x^i - \hat{\mu})^2 \tag{19}$$

$$Bias(\hat{\sigma}^2) = E[\hat{\sigma}^2 - \sigma^2] = \frac{n-1}{n}\sigma^2 - \sigma^2 = -\frac{1}{n}\sigma^2 \quad \hat{\sigma}^2 \text{ is biased} \tag{20}$$

$$\hat{\sigma}^2 \sim \frac{\sigma^2}{n-1}\chi^2_{n-1} \tag{21}$$

$$Var(\hat{\sigma}^2) = 2\sigma^2 \tag{22}$$

# Bias and Variance in Supervised Learning/Prediction

Similarities

- We use $\mathcal{D}_n$ to estimate $\hat{f}_n \in \mathcal{F}$
- $\mathcal{D}_n$ is random, hence so if $\hat{f}_n$.
- Bias and variance are properties of $\mathcal{F}$, and depend on $n$

  Exercise Consider linear regression $f(x) = \beta^T x$, with $N(0, \sigma^2)$ noise. What are the bias and variance of this predictor?

Differences

1. $\hat{f}$ is a function
2. We are interested in the predictions and not the parameters of $\hat{f}$.
3. We don't always assume $f^{true}$ exists.

- Several proposals to define bias and variance exist.
- What we need to know in this course/usually is qualitative

# Two definitions for bias in ML

1. Assuming $f^{true}$ exists
   - "Classical" framework
   - Typical example: Least Squares loss

2. (No assumption of $f^{true}$) Bias is **(in)ability of** $\mathcal{F}$ to fit the training set $\mathcal{D}_n$ (i.e. to make $\hat{L} = 0$)

3. In both cases, Variance is the variance of predictor $f(x)$ averaged over X

# The Bias-Variance decomposition for $L_{LS}$

▶ Assume true model $P_{Y|X}$

$$y = f(x) + \epsilon \quad \text{with} \quad \epsilon \sim iid, E[\epsilon] = 0, \quad Var(\epsilon) = \sigma^2, \; y, f(x), \epsilon \in \mathbb{R} \qquad (23)$$

▶ $\hat{f}_n$ is estimated from $\mathcal{D}_n$

▶ at $x$:

$$
\begin{aligned}
\text{MSE}(x) &= E_{P_{XY}^n}\left[\left(\hat{f}_n(x) - f^{true}(x)\right)^2\right] && (24) \\
&= E_{P_{XY}^n}\left[\left(\hat{f}_n(x) - E_{P_{XY}^n}[\hat{f}_n(x)] + E_{P_{XY}^n}[\hat{f}_n(x)] - f^{true}(x)\right)^2\right] && (25) \\
&= \underbrace{E_{P_{XY}^n}\left[\left(\hat{f}_n(x) - E_{P_{XY}^n}[\hat{f}_n(x)]\right)^2\right]}_{Var\hat{f}_n(x)} + \underbrace{E_{P_{XY}^n}}_{\text{deterministic}} \underbrace{\left[\left(E_{P_{XY}^n}[\hat{f}_n(x)] - f^{true}(x)\right)^2\right]}_{Bias^2(\hat{f}_n(x))} && (26) \\
&\quad + \underbrace{E_{P_{XY}^n}\left[\left(E_{P_{XY}^n}[\hat{f}_n(x)] - f^{true}(x)\right)\left(\hat{f}_n(x) - E_{P_{XY}^n}[\hat{f}_n(x)]\right)\right]}_{=0} && (27)
\end{aligned}
$$

▶ Note that $MSE(x) = E_{P_{Y|X}}[L_{LS}(y, \hat{f}_n(x))]$ <span style="color:purple">Exercise</span> Prove this

▶ Integrating over all $x \in \mathbb{R}$ w.r.t $P_X$

$$E_{P_X}[MSE(x)] \equiv E_{P_X}[E_{P_{Y|X}}[L_{LS}(y, \hat{f}_n(x))]] = E_{P_X}[Var(\hat{f}_n(x))] + E_{P_X}[Bias^2(\hat{f}_n(x))] \quad (28)$$

▶ Bias-Variance Decomposition of $L_{LS}$: $\boxed{L_{LS} = Var + Bias^2}$

# Case 2: Bias as model (mis)fit

- If no $f^{true}$ assumed, **bias** measures the (in)ability of the model class $\mathcal{F}$ to fit the data $\mathcal{D}_n$.
- Better fit $\Leftrightarrow$ less bias
- We measure the fit by the loss $L$ associated with the task, i.e $\hat{L}(\hat{f}_{\mathcal{D}_n}, \mathcal{D}_n)$
- **Bias**$(\mathcal{F}) = E_{P(X,Y)^n}[\hat{L}(\hat{f}_{\mathcal{D}_n}, \mathcal{D}_n)]$ (hence, bias is expected empirical loss).

  training error

- Richer or more complex models classes have less bias
  Example Bias( Linear ) > Bias( Quadratic )
  Example Bias( 1-NN ) < Bias( K-NN ), for $K > 1$

  $\mathcal{F}_1 \subset \mathcal{F}_2$
  $\{linear\} \subset \{quadratic\}$

- 
  Example Bias( Linear ) ? Bias( K-NN ) – depends on $P_{XY}$!

$$Eval \quad \mathcal{F} \quad by \quad \min_{f \in \mathcal{F}} \hat{L}(f, \mathcal{D})$$

$$\hat{L} = \frac{1}{n} \sum_{i=1}^{n} L(y^i, f(x^i))$$

---

[3] Not trivial, to find a reference.

# Case 2: Bias as model (mis)fit

- If no $f^{true}$ assumed, **bias** measures the (in)ability of the model class $\mathcal{F}$ to fit the data $\mathcal{D}_n$.
- Better fit $\Leftrightarrow$ less bias
- We measure the fit by the loss $L$ associated with the task, i.e $\hat{L}(\hat{f}_{\mathcal{D}_n}, \mathcal{D}_n)$
- **Bias**$(\mathcal{F})= E_{P(X,Y)^n}[\hat{L}(\hat{f}_{\mathcal{D}_n}, \mathcal{D}_n)]$ (hence, bias is expected empirical loss).

- Richer or more complex models classes have less bias
  Example  Bias( Linear ) $>$ Bias( Quadratic )
  Example  Bias( 1-NN ) $<$ Bias( K-NN ), for $K > 1$
-
  Example  Bias( Linear ) ? Bias( K-NN ) – depends on $P_{XY}$!
- In modern ML we consider sequences of model classes that can be ordered
    - by inclusion
$$\mathcal{F} \subset \mathcal{F}' \quad \text{then bias}(\mathcal{F}) \geq \text{bias}(\mathcal{F}') \tag{29}$$
      Example Linear $\subset$ Quadratic, . . . CART( $L$ leaves ) $\subset$ CART( $L+1$ leaves) . . . , Neural net ($L$ layers) $\subset$ . . .
    - by complexity
$$\text{complexity}(\mathcal{F}) \leq \text{complexity}(\mathcal{F}') \quad \text{then bias}(\mathcal{F}) \geq \text{bias}(\mathcal{F}') \tag{30}$$
      Example . . . complexity(Kernel($h$)) $\nearrow$ for $h \downarrow$, complexity(Ridge Regression, Lasso( $\lambda$ )) $\uparrow$ for $\lambda \downarrow$, complexity(Linear with margin $R$) $\uparrow$ for $R \downarrow$
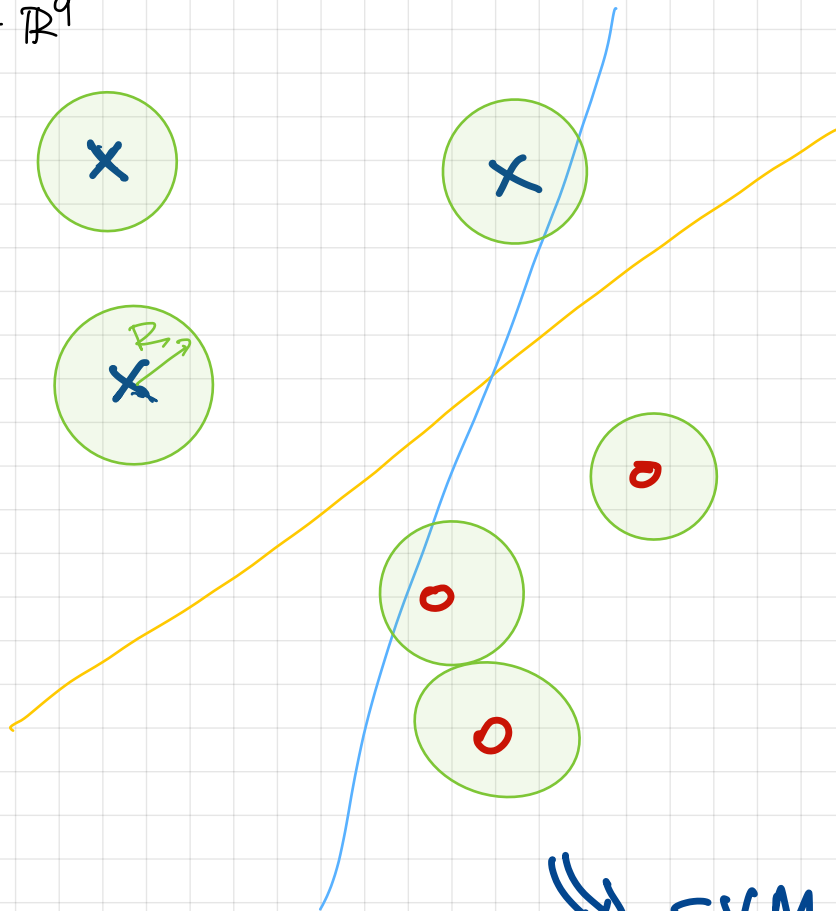- Larger data are harder to fit (hence more bias on average)[3]

*When $n \nearrow$*
*$\mathcal{F}$ the same*
*Bias-Var tradeoff changes*

---

[3] Not trivial, to find a reference.

$X \in \mathbb{R}^d$

$$\text{VC dim} \overbrace{(\text{Linear})}^{d_0} = d + 1$$

$$\alpha_R = \{ \text{linear, margin } R \}$$

$$\alpha_{R_1} < \alpha_{R_2} < \alpha_0$$

$$R_1 > R_2$$

$\rightsquigarrow$ SVM

# Sampling variance

▶ Intuition: if we draw two different data sets $\mathcal{D}, \mathcal{D}' \sim P_{XY}$ (from the same distribution) we will obtain different predictors $f, f'$. Variance measures how different the predictions of $f, f'$ can be on average.

▶ **Variance** at $x = Var_{P_{XY}^n}(\hat{f}_{\mathcal{D}_n}(x))$, where the randomness is over the sample $\mathcal{D}_n$

▶ **Variance** associated with predictor class $\mathcal{F}$ is the expectation over $P_X$ of the variance at $x$, i.e $E_{P_X}[Var_{P_{XY}^n}(\hat{f}_{\mathcal{D}_n}(x))]$

▶ Variance depends on $n$, $\mathcal{F}$, and the data distribution $P_{XY}$  Exercise If $P_{Y|X}$ is deterministic for all $x$, does it mean that the variance is 0?

▶ Richer model classes are subject to more variance

$$\mathcal{F} \subset \mathcal{F}' \quad \text{then } Var(\mathcal{F}) \leq Var(\mathcal{F}') \quad \text{for any } f^*$$

# Variance, bias and model complexity

▶ Synonyms: rich class = complex model = flexible model = high modeling power = many degrees of freedom = many parameters

▶ Evaluating the model complexity[4]/number of free parameters of a model class $\mathcal{F}$ is usually a difficult problem!

Non-parametric models # parameters depends on $P_{XY}$, smoothing parameter and $n$
Parametric models # parameters NOT always equal to the number of parameters of $f$!

Example the classifier $f(x) = \mathrm{sgn}(\alpha x), x, \alpha \in \mathbb{R}$ depends on one parameter $\alpha$ but has $\infty$ degrees of freedom[5]!

Example the linear classifier and regressor on $\mathbb{R}^d$ has (no more than) $n + 1$ degrees of freedom

Example the complexity of a two layer neural net with $m$ fixed is not known (but there are approximation results); the number of weights in $f$ is obviously $(m + 1)(d + 1) + 1$

Example For K-NN, the variance increases when $K$ decreases

Example For pruned Decision Tree, the variance increases whith the number of levels

▶ The variance of a predictor increases with the complexity of $\mathcal{F}$.
▶ But complexity is the opposite of bias, so bias decrease with the complexity of $\mathcal{F}$
▶ This is known as the **Bias-Variance tradeoff**

---

[4]There are several definitions of model complexity, but this holds for all definitions I know
[5]See VC-dimension later

# The Bias-Variance tradeoff   *Case II*

| Wanted property (for an $\mathcal{F}$) | unwanted consequence of $\mathcal{F}$ not satisfying this property | what to do |
|---|---|---|
| to fit $\mathcal{D}$ well | Bias       *overas overfitting* | increase complexity |
| to be robust to sampling noise | Variance   *under fitting* | decrease complexity |

The **bias-variance tradeoff** is the observation that the better a predictor class $\mathcal{F}$ is able to fit any given sample, the more sensitive the selected $f$ will be to sampling noise.
In this course we will learn some ways of balancing these desired properties (or these undesired consequences).

$$n \uparrow \quad \text{always better !!}$$

$$L \;\leq\; \hat{L} \;+\; O\!\left(\tfrac{1}{\sqrt{n}}\right) \quad \text{w.p. } 1-\delta$$

complexity of $\overline{\mathcal{F}}$

Expected Loss     Bias     Variance

## Examples, examples. . .

### Example (*K*-nearest neighbor classifiers)

The 1-NN can fit any data set perfectly (every data point is it's own nearest neighbor). But for $K > 1$, the $K$-NN may not be able to reproduce any pattern of $\pm 1$ in the labels. Hence its bias is larger than the bias of the 1-NN classifier. With the variance, the opposite happens: as $K$ the number of neighbors increases, the decision regions of the $K$-NN classifier become more stable to the random sampling effects. Thus, the variance decreases with $K$.

### Example (Linear vs quadratic vs cubic . . . predictors)

The quadratic functions include all linear functions, the cubics include all quadratics, and so on. Linear classifiers will have more bias (less flexibility) than quadratic classifiers. On the other hand, the variance of the linear classifier will be lower than that of the quadratic. The case of regression is even more straightforward: if we fit the data with a higher degree polynomial, the fit will be more accurate, but the variation of the polynomial $f(x)$ for $x$ values not in the training set will be higher too.

### Example (Kernel regression)

## Examples, examples... (2)

The bias-variance tradeoff can be observed on a continuous range for **kernel regression**. When the kernel width $h$ is near 0, $f(x)$ from Lecture 1, equation (25) will fit the data in the training set exactly [Exercise: prove this], but will have high variance. When $h$ is large, $f(x^i)$ will be smoothed between $x^i$ and the other data points nearby, so it may be some distance from $y^i$. However, precisely because $f(x)$ is supported by a larger neighborhood, it will have low variance. [Exercise: find some intuitive explanations for why this is true] Hence, the smoothness parameter $h$ controls the trade-off between bias and variance.
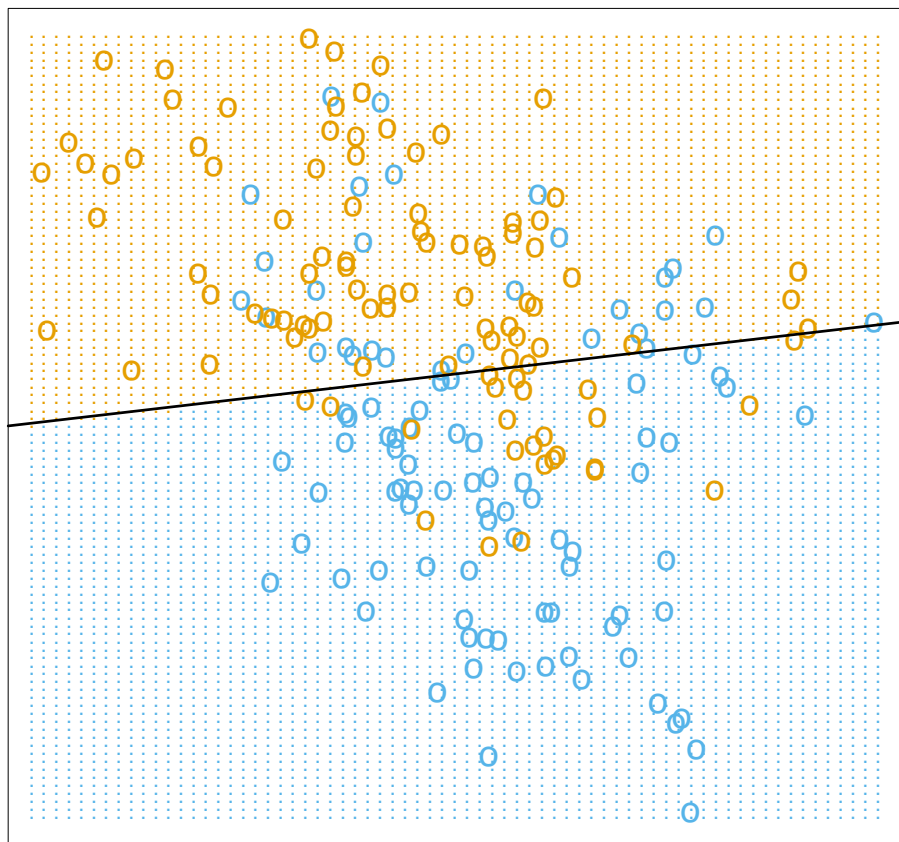
## Example (Regularization)

The same can be observed if one considers equation (15). For $\lambda = 0$, one choses $f$ that best fits the data (minimizes $\hat{L}$. For $\lambda \to \infty$, $f$ is chosen to minimize the penalty $J$, disregarding the data completely. The latter case has 0 variance, but very large bias. Between these extreme cases, the parameter $\lambda$ controls the amount in which we balance fitting the data (variance) with pulling $f$ towards an a-priori "good" (bias).

# Overfitting and Underfitting

▶ Bias and variance are properties of the model class $\mathcal{F}$ (sometimes toghether with the learning algorithm – more about this later). They are not properties of the parameters of $f$ (e.g $\beta$), and not of a particular $f \in \mathcal{F}$.

▶ Variance decreases to 0 with $n$, but bias may not. This implies that for larger sample sizes $n$, the trade-off between variance and bias changes, and typically the "best" trade-off, aka the best model, will have larger complexity.

▶ **Overfitting**= is the situation of small bias and too much variance (i.e. $\mathcal{F}$ is too complex). In practice, if a learned predictor $f$ has low $\hat{L}(f)$ but significantly higher $L(f)$, we say that the model has *overfit* the data $\mathcal{D}$. (Of course we cannot know $L(f)$ directly, and a significant amount of work in statistics is dedicated to predicting $L(f)$ for the purpose of chosing the best model.)

▶ **Underfitting**=bias is too high, or the model is too simple (a.k.a has too few degrees of freedom). [Exercise: what do you expect to see w.r.t. $\hat{L}(f)$ v.s. $L(f)$ for an underfitted model?]
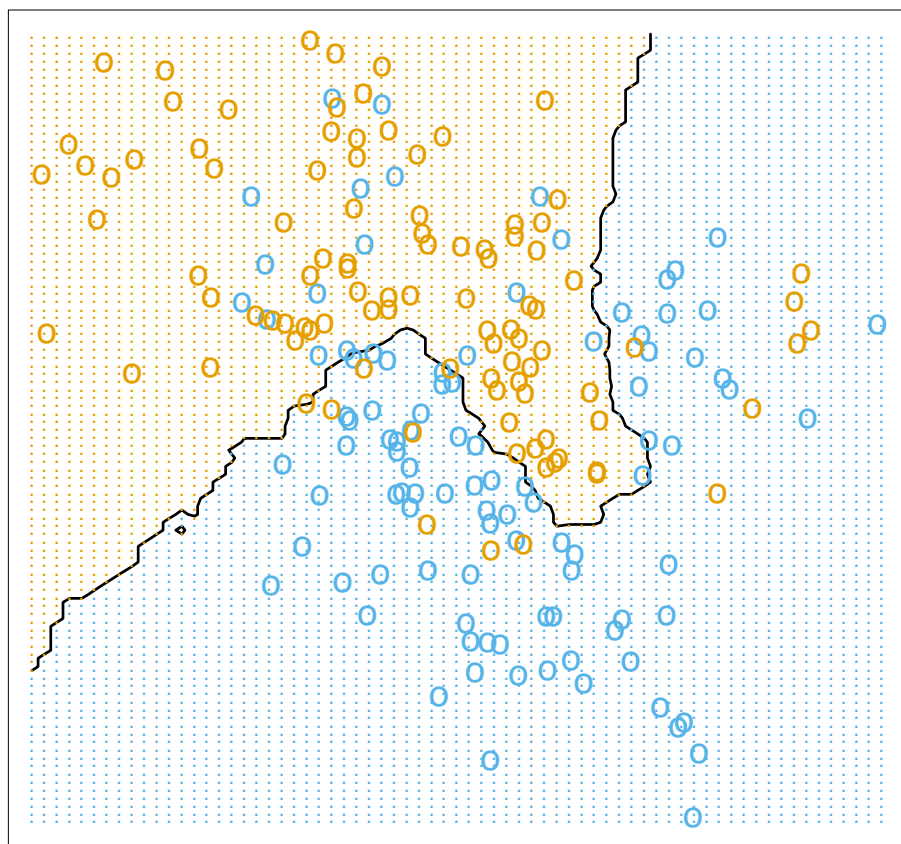
Complexity, even though there are variations in its definition, and although it is not known exactly for most model classes, is at the core of learning theory, the part of statistical theory that gives provable results about the expected loss of a predictor.
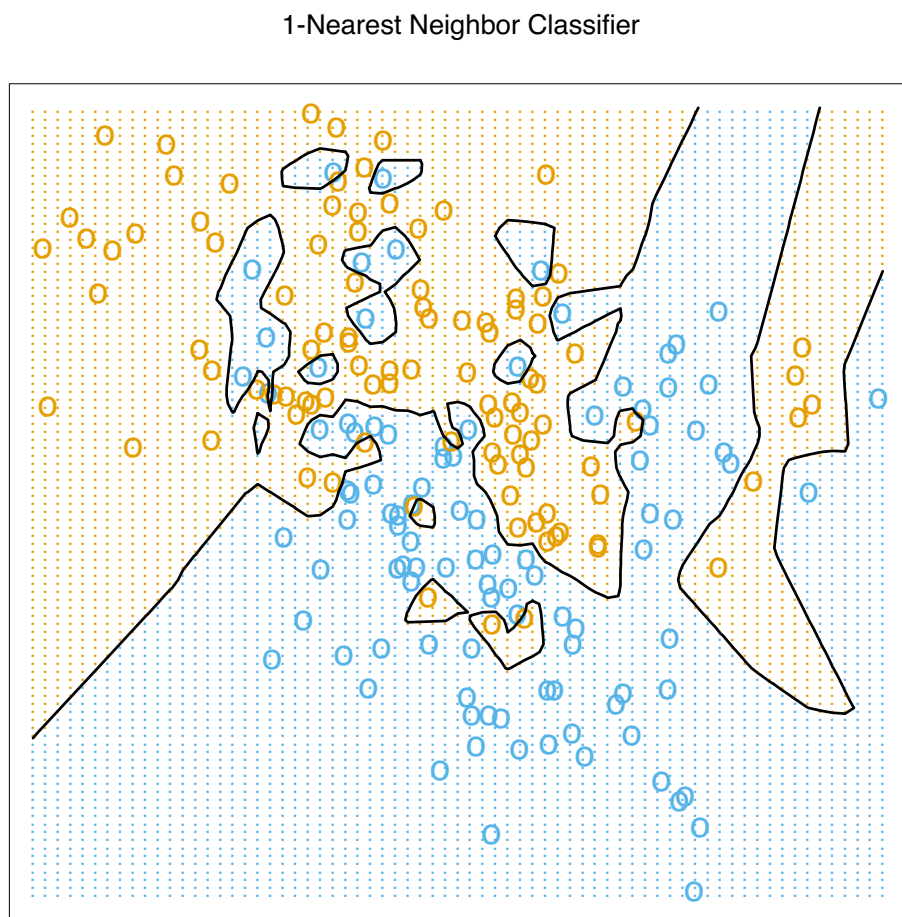
Linear Regression of 0/1 Response



**FIGURE 2.1.** *A classification example in two dimensions. The classes are coded as a binary variable (BLUE = 0, ORANGE = 1), and then fit by linear regression. The line is the decision boundary defined by $x^T \hat{\beta} = 0.5$. The orange shaded region denotes that part of input space classified as ORANGE, while the blue region is classified as BLUE.*
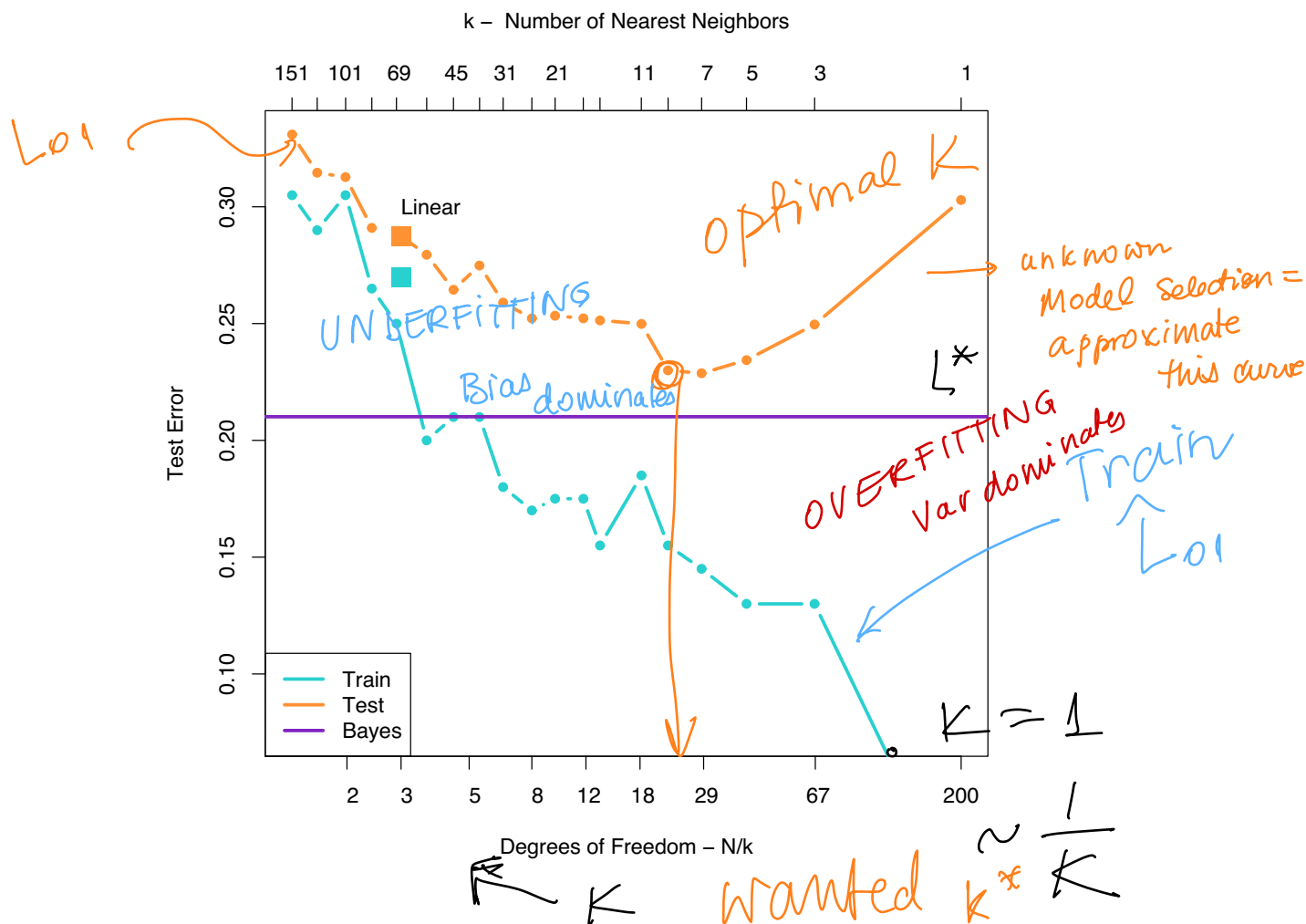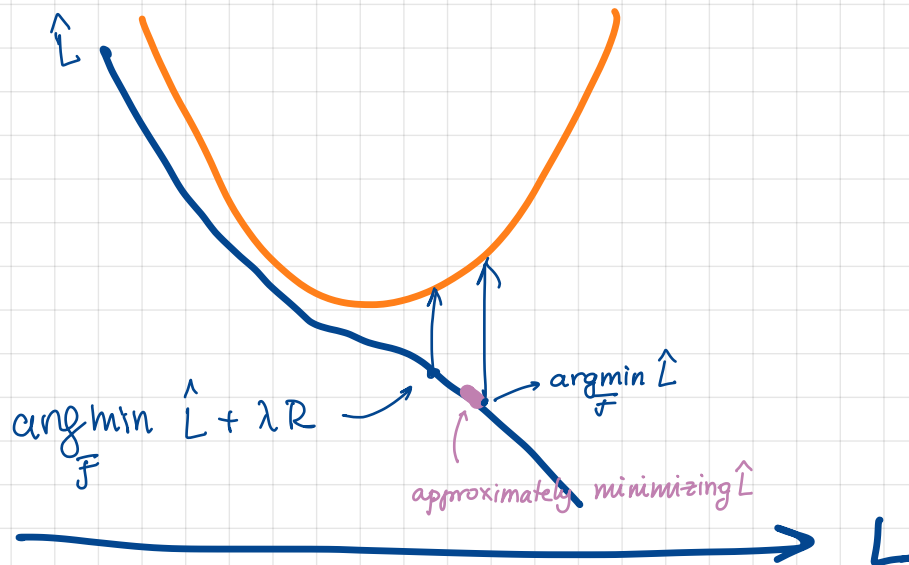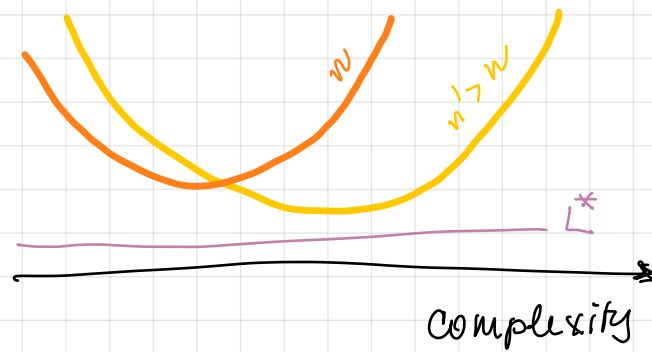
15-Nearest Neighbor Classifier



**FIGURE 2.2.** *The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable* ($\text{BLUE} = 0, \text{ORANGE} = 1$) *and then fit by 15-nearest-neighbor averaging as in (2.8). The predicted class is hence chosen by majority vote amongst the 15-nearest neighbors.*
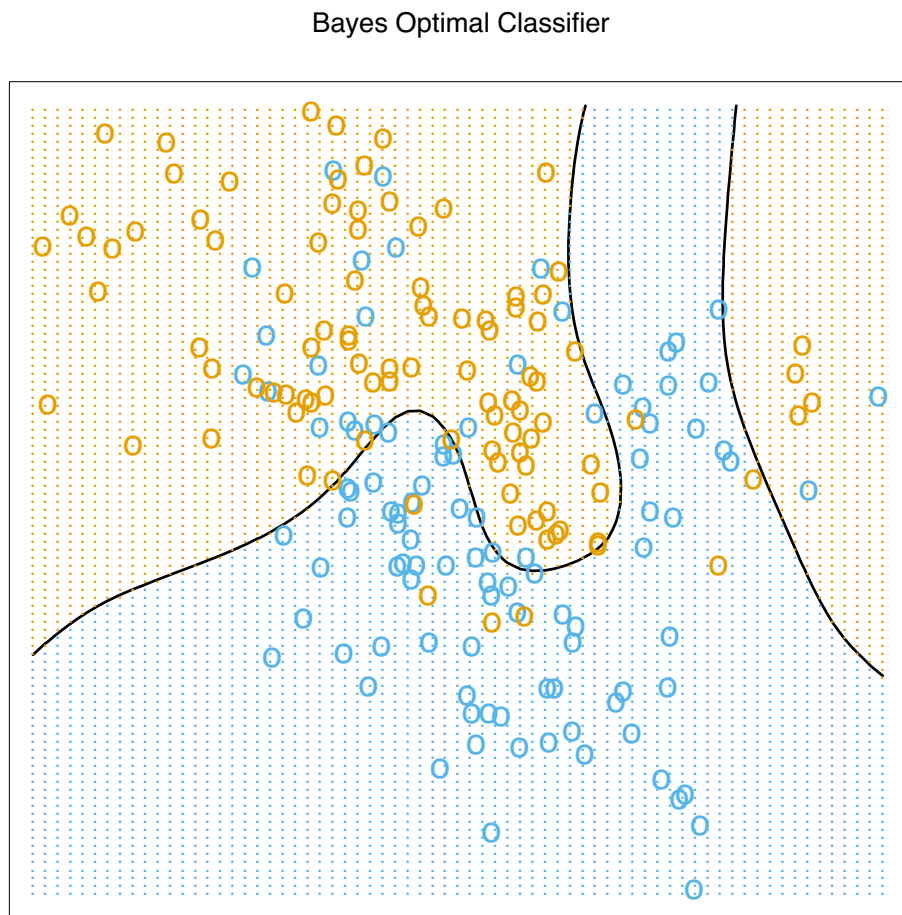
1-Nearest Neighbor Classifier



**FIGURE 2.3.** *The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable (BLUE = 0, ORANGE = 1), and then predicted by 1-nearest-neighbor classification.*
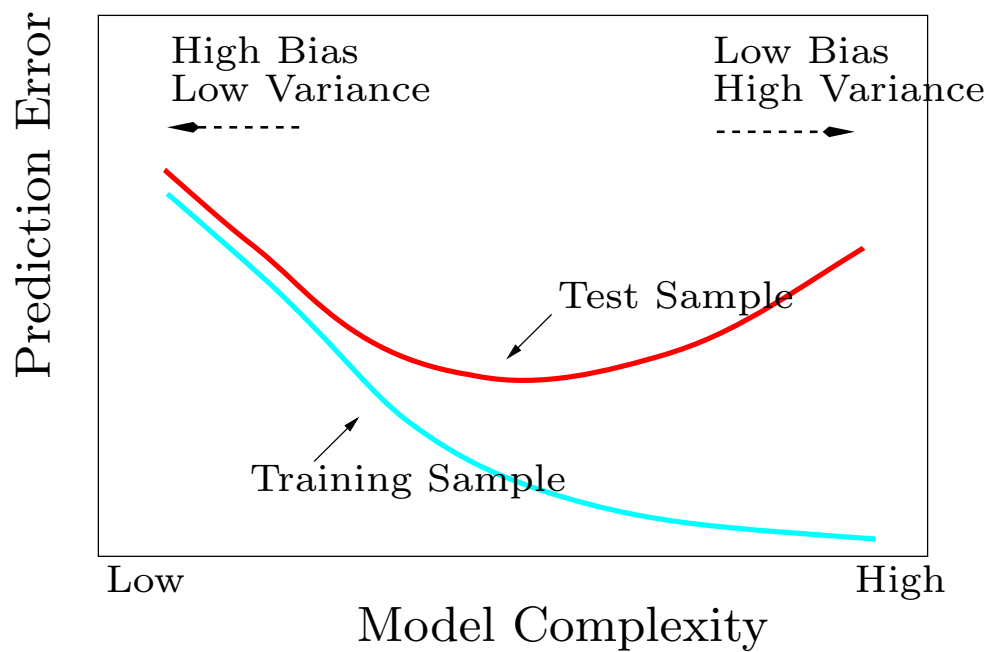
**FIGURE 2.4.** *Misclassification curves for the simulation example used in Figures 2.1, 2.2 and 2.3. A single training sample of size* 200 *was used, and a test sample of size* 10,000. *The orange curves are test and the blue are training error for k-nearest-neighbor classification. The results for linear regression are the bigger orange and blue squares at three degrees of freedom. The purple line is the optimal Bayes error rate.*

$n$

$n' > n$

$L^*$

complexity

$\hat{L}$

argmin $\hat{L} + \lambda R$
$F$

argmin $\hat{L}$
$F$

approximately minimizing $\hat{L}$
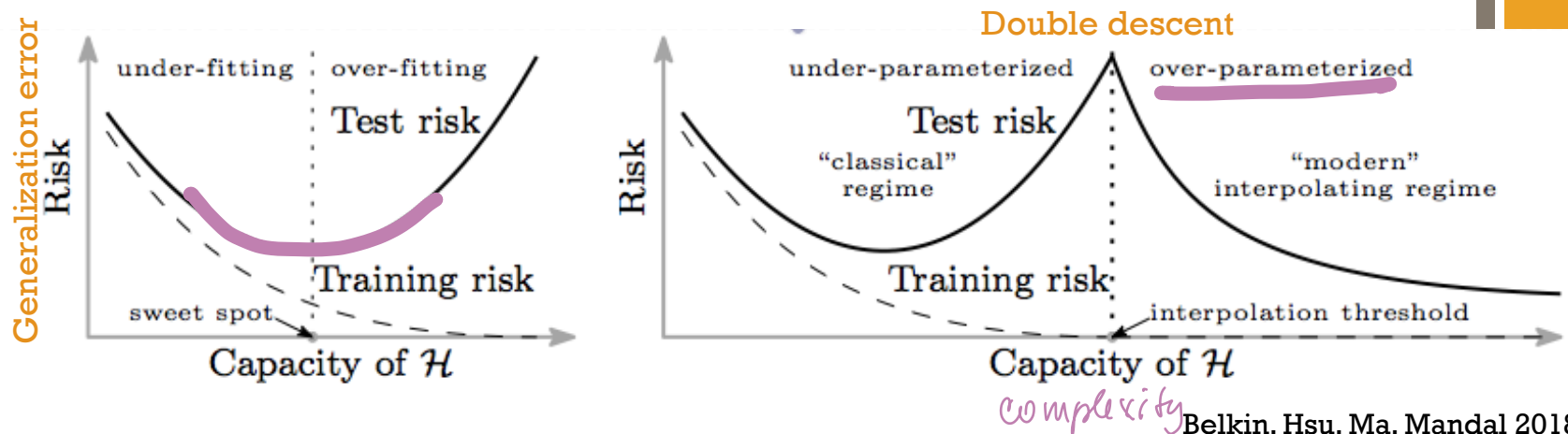
$L$

Bayes Optimal Classifier



**FIGURE 2.5.** *The optimal Bayes decision boundary for the simulation example of Figures 2.1, 2.2 and 2.3. Since the generating density is known for each class, this boundary can be calculated exactly (Exercise 2.2).*

**FIGURE 2.11.** *Test and training error as a function of model complexity.*

# What is observed



Belkin, Hsu, Ma, Mandal 2018

- Classical regime p < N

- Modern/Deep Learning/High dimensional regime N > n
    - Think N fixed, p increases, gamma=p/N
    - Training error = 0 (interpolation)
    - Test error decreases with p (or gamma)

# Lecture Notes III – Neural Networks

Marina Meilă
mmp@stat.washington.edu

Department of Statistics
University of Washington

October, 2023

Two-layer Neural Networks    *Logistic Regression*

Multi-layer neural networks    *Backpropagation*

A zoo of multilayer networks    *Architectures*

**Reading** HTF Ch.: 11.3 Neural networks, Murphy Ch.: (16.5 neural nets), Bach Ch.: –, Deep Learning Book (Goodfellow, Bengio, Courville) 6.1-4, ResNet 7.6, ConvNet 9., Autoencoders 14.1, Dive Into Deep Learning 4.1-4.3.

# Two-layer Neural Networks

▶ The **activation function** (a term borrowed from neuroscience) is any continuous, bounded and strictly increasing function on $\mathbb{R}$. Almost universally, the activation function is the **logistic** (or **sigmoid**)

$$\phi(u) = \frac{1}{1 + e^{-u}} \tag{1}$$

because of its nice additional computational and statistical properties.

▶ We build a **two-layer neural network** in the following way:

| | | |
|---|---|---|
| Inputs | $x_k$ | $k = 1 : d$  →"layer" 0 |
| Bottom layer[1] | $z_j = \phi(w_j^T x)$ | $j = 1 : m, \ w_j \in \mathbb{R}^d$ →"layer" 1 (hidden) |
| Top layer | $f = \beta^T z$ | $\beta \in \mathbb{R}^m$ → "layer" 2 (output) |
| Output | $f$ | $\in [0, 1]$ |

$m$ = size of "hidden" layer

In other words, the neural network implements the function

$$\textbf{output } f(x) = \sum_{j=1}^{m} \beta_j z_j = \sum_{j=1}^{m} \beta_j \phi(\underbrace{\sum_{k=1}^{d} w_{kj} x_k}_{z_j}) \ \in (-\infty, \infty) \tag{2}$$

Note that this is just a linear combination of logistic functions.

$$x \in \mathbb{R}^d \qquad z_{1:m} = \phi(w_{1:m}^T x) = \textit{logistic regression}$$
$$w_{1:m} \in \mathbb{R}^d \qquad \beta \in \mathbb{R}^m$$

---

[1] In neural net terminology, each variable $z_j$ is a **unit**, the bottom layer is **hidden**, while top one is **visible**, and the units in this layer are called hidden/visible units as well. Sometimes the inputs are called **input units**; imagine neurons or individual circuits in place of each $x, y, z$ variable.