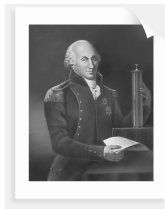
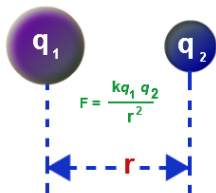


Introduction to Machine Learning

Yet Looking for Tools to Model

- We have all seen simple physics equations, discovered hundreds of years ago
- Since then we have been curious about how systems work and to predict their behavior before constructing them



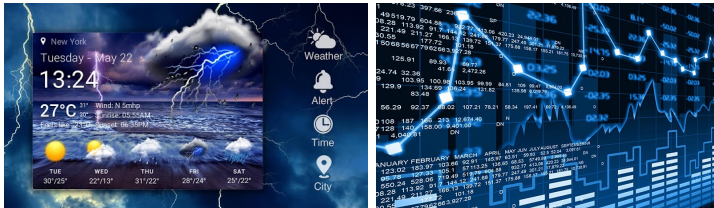
- Back then, many derivations were based on experiments
- For instance deriving an equation like $F = 8.99 \times 10^9 \frac{q_1 q_2}{r^2}$ could involve fixing all the parameters except one and evaluating the response behavior
- In this simple example q_1, q_2, r are variables and F is a response

Yet Looking for Tools to Model

- What has changed?
 - Back then, we did not have accurate measurement systems but now we have very accurate measurement sensors for pretty much all physical phenomena
 - Back then, we were only able to collect limited data, now we are able to collect huge data
 - Back then, we did not have computational resources, now we have supercomputers
 - Back then, we always liked more compact models for an easier application, now we can have super-complex models and ask the computers to apply them
- **Then why not model large and complex systems??!**

Yet Looking for Tools to Model

- Models that can predict the weather, the stock market value, political interactions, ...



- In our century, having more data means more power to predict!

More Examples: Spam Detection

- data from 4601 emails sent to an individual (named George, at HP labs, before 2000). Each is labeled as spam or email
- goal: build a customized spam filter
- input features: relative frequencies of 57 of the most commonly occurring words and punctuation marks in these email messages

	george	you	hp	free	!	edu	remove
spam	0.00	2.26	0.02	0.52	0.51	0.01	0.28
email	1.27	1.27	0.90	0.07	0.11	0.29	0.01

Average percentage of words or characters in an email message equal to the indicated word or character. We have chosen the words and characters showing the largest difference between **spam** and **email**

More Examples: Optical Character Recognition

identify hand-written characters based on the combination of all pixel values

label = 5



label = 0



label = 4



label = 1



label = 9



label = 2



label = 1



label = 3



label = 1



label = 4



label = 3



label = 5



label = 3



label = 6



label = 1



label = 7



label = 2



label = 8



label = 6

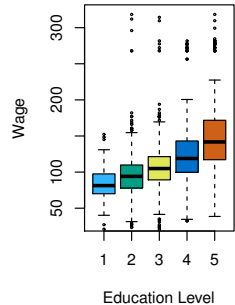
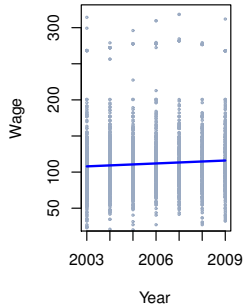
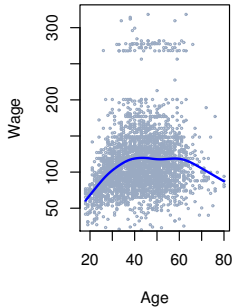


label = 9



More Examples: Salary Prediction

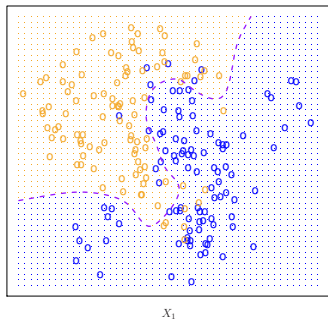
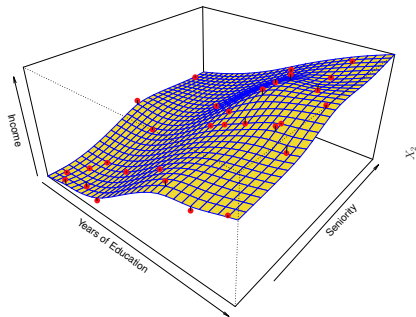
- Establish the relationship between salary and demographic variables in population survey data



Supervised Learning

- Recall the equation $F = 8.99 \times 10^9 \frac{q_1 q_2}{r^2}$ or more generally $Y = f(\mathbf{X})$ where $\mathbf{X} = (X_1, \dots, X_p)^\top$
 - In the case above, $Y = F$ and $\mathbf{X} = (q_1, q_2, r)^\top$ *features*
 - Outcome measurement Y (also called dependent variable, response, target)
 - Vector of p predictor measurements \mathbf{X} (also called inputs, regressors, covariates, features, independent variables)
 - In the **regression** problem, Y is quantitative (e.g price, blood pressure)
 - In the **classification** problem, Y takes values in a finite, unordered set (survived/died, digit 0-9, cancer class of tissue sample) *Discrete Number*
 - We have training data $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_N, Y_N)$. These are observations (examples, instances) of these measurements
- [**Bold** faces correspond to vectors]

Examples of Regression and Classification



On the basis of the training data we would like to:

- Accurately predict unseen test cases
 - Understand which inputs affect the outcome, and how
 - Assess the quality of our predictions and inferences
- Normally from the available examples $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_N, Y_N)$ we choose N_1 of them (around 80%) for training and $N - N_1$ of them for the later evaluation of the model (testing)

< 80% training dataset
20% evaluation dataset

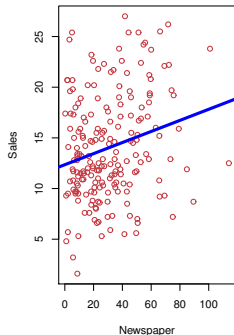
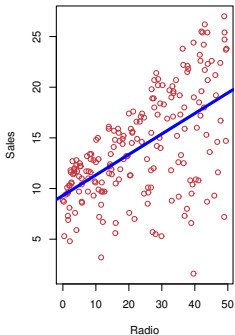
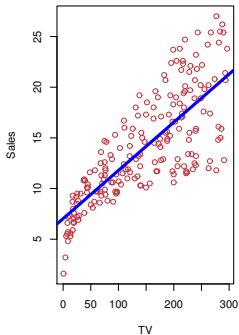
Unsupervised Learning

- No outcome variable, just a set of predictors (features) measured on a set of samples.
- objective is more fuzzy — find groups of samples that behave similarly, find features that behave similarly, find linear combinations of features with the most variation.
- difficult to know how well you are doing.
- different from supervised learning, but can be useful as a pre-processing step for supervised learning.
- Examples
 - Cocktail Party problem [e.g., see this <https://tinyurl.com/33sma6yt>]
 - [Just like the Facebook tagging system] you have collection of photos of multiple people, without information who is on which. You want to divide this dataset into multiple piles, each with photos of one individual
 - Generative Adversarial Networks (GANs)

More on the Basics of Statistical Learning

An Example

Shown are Sales vs the advertising budget on TV, Radio and Newspaper.
Blue linear-regression line fits separately to each:



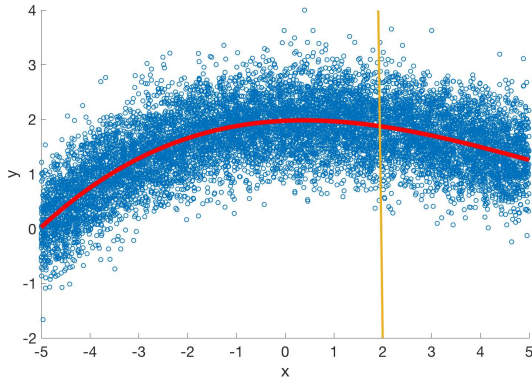
Can we predict Sales using these three? Basically,
 $\text{sales} = f(\text{TV}, \text{Radio}, \text{Newspaper})$

Why do we Need Predictive Models

- With a good fit we can make predictions of Y at new points $\mathbf{X} = \mathbf{x}$
- We can understand which features (i.e., components of $\mathbf{X} = (X_1, X_2, \dots, X_p)^\top$) are important in explaining Y
- For instance in the sales example TV, Radio and Newspaper can be important but “which restaurant is close to the company” might be an irrelevant feature
- Depending on the complexity of the fit, we may be able to understand how each component X_j of \mathbf{X} affects Y

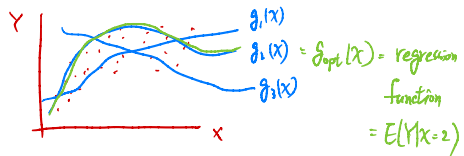
Best Predictive Model if We Had Enough Data

- Suppose we had many data samples (X, Y) as below
- What would have been a good estimate of Y for $X = 2$



- In mathematical terms what we are after is $E(Y|X = 2)$

Best Predictive Model if We Had Enough Data



- In an ideal setting of having many data, the reasonable value of Y corresponding to $X = \mathbf{x}$ is the average of all responses at $\mathbf{X} = \mathbf{x}$
- It can be mathematically shown that $f(\mathbf{x}) = E(Y|\mathbf{X} = \mathbf{x})$ is the function that minimizes $E[(Y - g(\mathbf{X}))^2|\mathbf{X} = \mathbf{x}]$ over all functions g at all points $\mathbf{X} = \mathbf{x}$.
- This ideal $f(\mathbf{x}) = E(Y|\mathbf{X} = \mathbf{x})$ is called the **regression function**
aka "Good's Model"

Question: Can anyone show that $f(x) = E(Y|X = x)$ is the function that minimizes $E[(Y - g(X))^2|X = x]$ over all functions g ?

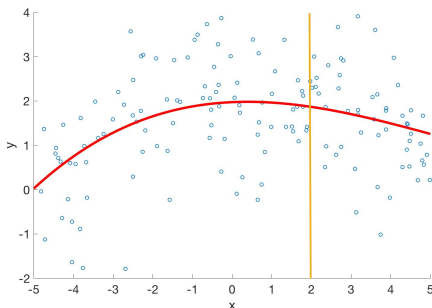
$$C = E_{|X} [(y - g(x))^2] = E_{|X} (y^2 - 2y \cdot g + g^2) = E_{|X}(y^2) - 2g \cdot E_{|X}(y) + g^2$$

$$\frac{\partial C}{\partial g} = 0 - 2 E_{|X}(y) + 2g = 0 \Rightarrow g_{\text{opt}} = E_{|X} = E(y|X=x)$$

What Happens in Practice

In practice we don't have enough data to have access to f

Even if that could have been the case, doing it for all points \mathbf{x} in the space is not practical



One immediate suggestion is considering neighborhoods around \mathbf{x}

Nearest neighbor averaging can only be good for small p (say < 4) and large N . [Example:] Points $\mathbf{x} = (0.9, \dots, 0.9) \in \mathbb{R}^{100}$ and $\hat{\mathbf{x}} = (1, \dots, 1) \in \mathbb{R}^{100}$ are quite far

Fundamental Limit

- As we observed the best we can do is the regression function

$$Y = f(\mathbf{x}) + \epsilon$$

- $f(\mathbf{x})$ can be considered as the actual physical model that generates data at $\mathbf{X} = \mathbf{x}$ and ϵ as the noise and uncertainty
- There is nothing we can do about ϵ , since it is a random variable not under our control
- Even in practice we cannot estimate $f(\mathbf{x})$ for $\mathbf{X} = \mathbf{x}$ and our estimate is something like $\hat{f}(\mathbf{x})$ different than $f(\mathbf{x})$
- Hence, the overall prediction error at a point $\mathbf{X} = \mathbf{x}$ is

$$E \left((Y - \hat{f}(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x} \right) = \underbrace{\left(f(\mathbf{x}) - \hat{f}(\mathbf{x}) \right)^2}_{\text{GM}} + \text{Var}(\epsilon)$$

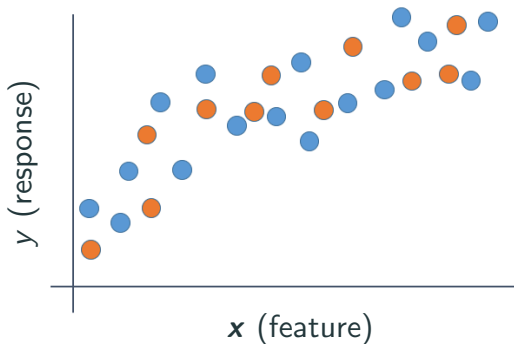
Actual Model

- First term can be improved by improving $\hat{f}(\mathbf{x})$, but the second term is constant and nothing can be done about it

$$\begin{aligned}
 E_{|X} [y - \hat{f}(x)]^2 &= E_{|X} (f + e - \hat{f})^2 = E_{|X} [(f - \hat{f})^2 + 2e(f - \hat{f}) + e^2] \\
 &= (f - \hat{f})^2 + 2(f - \hat{f}) \cdot E[e] + E(e^2) \\
 &= (f - \hat{f})^2 + \text{Var}(e)
 \end{aligned}$$

How to Assess Model Accuracy in Practice

- As we said before, we don't normally have access to $f(\mathbf{x})$, so how can we assess the accuracy of an estimated model $\hat{f}(\mathbf{x})$?
- Usually we split the data into a **training set** and a **test set** to later evaluate the accuracy



Parametric vs Nonparametric Models

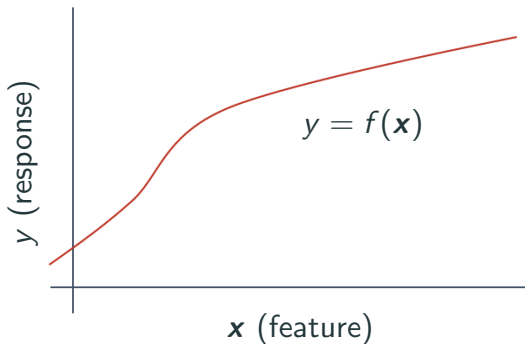
- In parametric models, we use some predefined forms for the model \hat{f} that we would like to fit
- For instance in linear regression we use the following general form and the fitting aims to find the coefficients

$$Y = f_L(\mathbf{X}) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots \beta_p X_p$$

- Non-parametric methods do not make explicit assumptions about the functional form of the fit. Instead they seek an estimate of f that gets as close to the data points as possible without being too rough or wiggly
- Thin-plate spline is an example of nonparametric fit, where you can control the smoothness of the fit vs its flexibility to match the training data

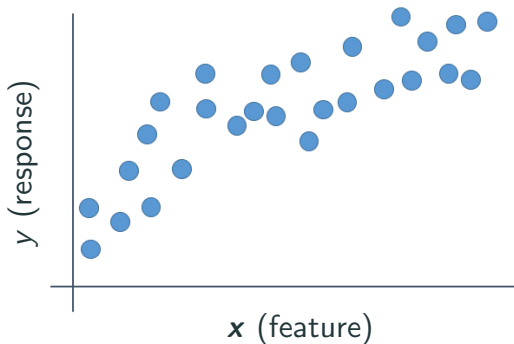
Basics of Model Fitting

- We have a (physical) system which produces a response y to an input x



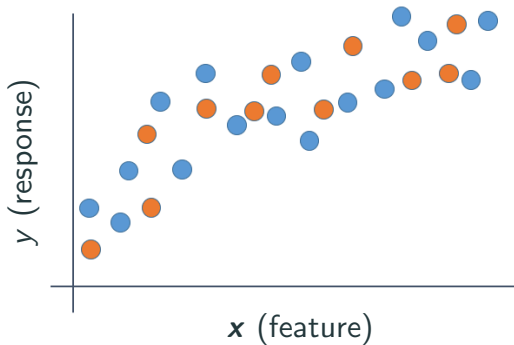
Basics of Model Fitting

- The system is a black box and all we have are samples of input/output



Basics of Model Fitting

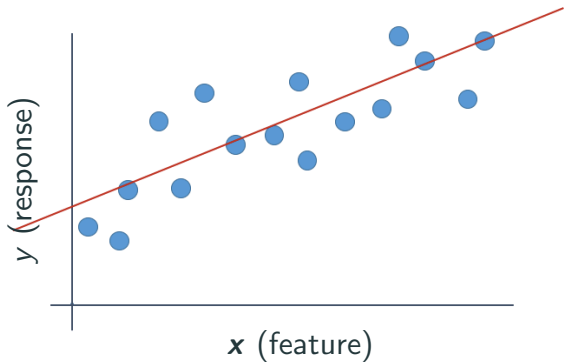
- We may try fitting a model to the data



- We split the data into a **training set** and a **test set** to later evaluate the accuracy

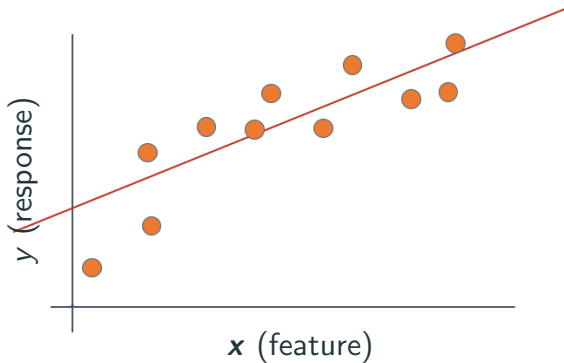
Basics of Model Fitting

- A linear model: $y = w_1x + w_0$



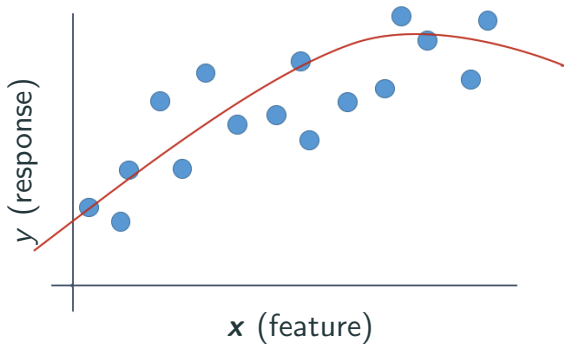
Basics of Model Fitting

- Testing the linear model: $y = \hat{f}(x) = w_1x + w_0$



Basics of Model Fitting

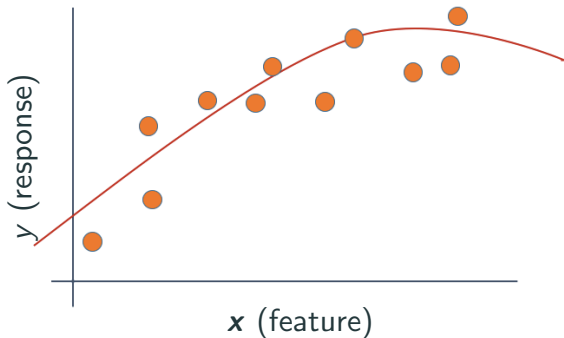
- Second order polynomial model: $y = \hat{f}(x) = w_2x^2 + w_1x + w_0$



Basics of Model Fitting

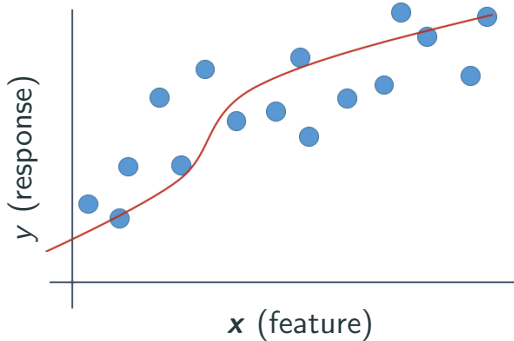
- Testing the second order polynomial model:

$$y = \hat{f}(x) = w_2x^2 + w_1x + w_0$$



Basics of Model Fitting

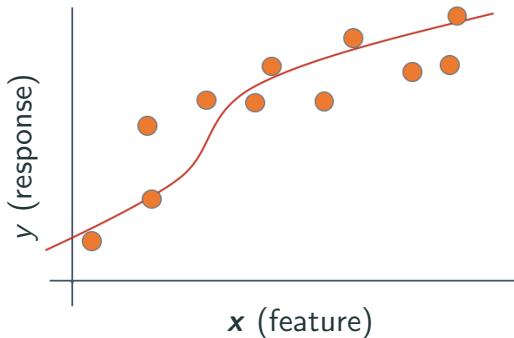
- Third order polynomial model: $y = \hat{f}(x) = w_3x^3 + w_2x^2 + w_1x + w_0$



Basics of Model Fitting

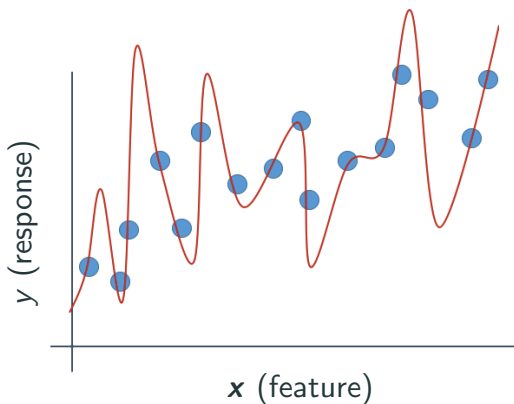
- Testing the third order polynomial model:

$$y = \hat{f}(x) = w_3x^3 + w_2x^2 + w_1x + w_0$$



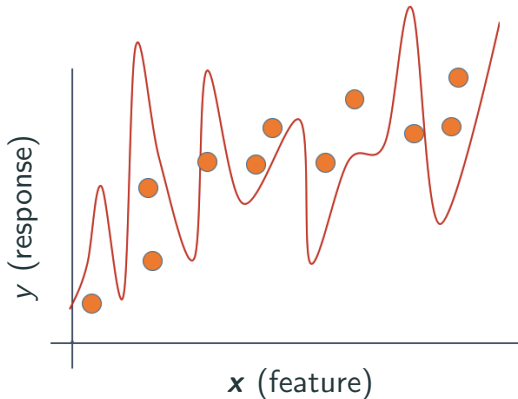
Basics of Model Fitting

- High order polynomial model: $y = \hat{f}(x) = \sum_{p=0}^P w_p x^p$ or a thin-plate spline with minor smoothness control



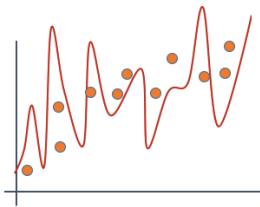
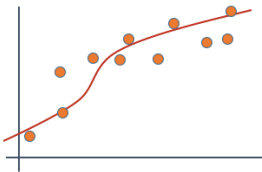
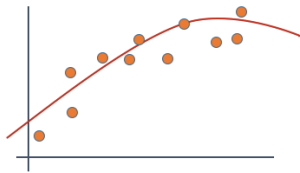
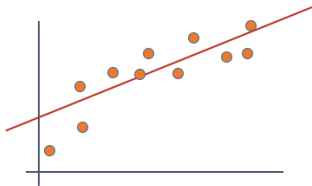
Basics of Model Fitting

- Testing the high order polynomial model: $y = \hat{f}(x) = \sum_{p=0}^P w_p x^p$ or a thin-plate spline with minor smoothness control



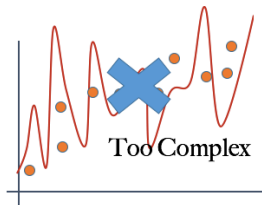
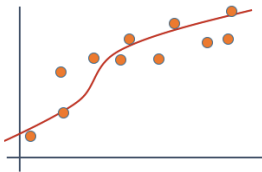
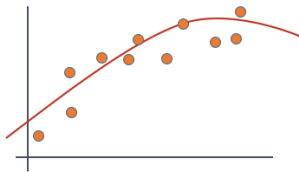
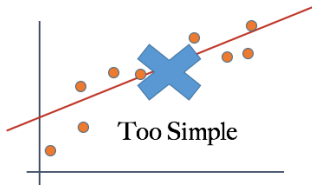
Basics of Model Fitting

- Which model to pick? Underfitting vs overfitting



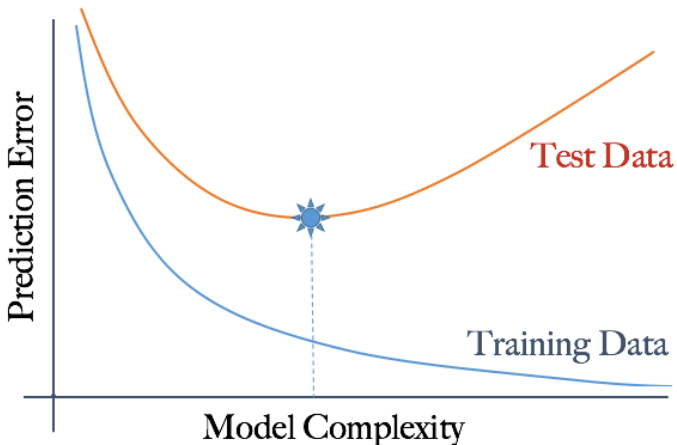
Basics of Model Fitting

- Which model to pick? Underfitting vs overfitting



Typical Curve

- The balance between model flexibility and test error typically looks as below



Assessing Model Accuracy

- Once we fit our model \hat{f} we want to evaluate its performance
- Of course the evaluation should not be with respect to the training data, since that would not be a fair evaluation
- Instead the evaluation is performed on the test data (that has not been used in the fitting)
- Mean Squared Error = $\text{Ave}_{i \in \text{Test}}(Y_i - \hat{f}(\mathbf{x}_i))^2$

Bias vs Variance(Theory)

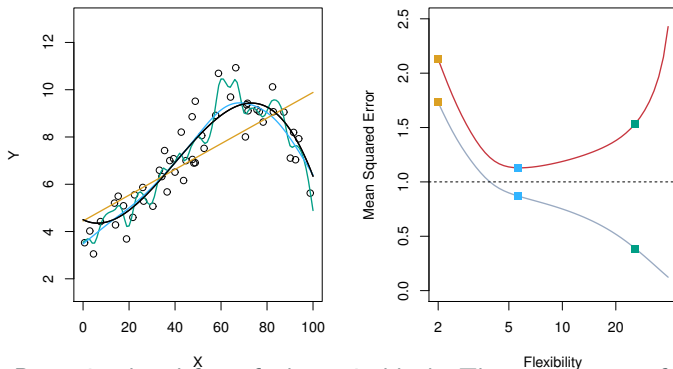
- Suppose we fit a model \hat{f} and we have a large test set to accurately calculate the MSR with respect to the test data
- Assuming that (\mathbf{x}_t, Y_t) is a fixed point in the test set and T is the training set based on which \hat{f} is derived, we generally have

$$E_{tot} \left(Y_t - \hat{f}(\mathbf{x}_t) \right)^2 = \overset{\text{Model Complexity}}{\text{var}(\hat{f}(\mathbf{x}_t))} + \overset{\text{over-fitting evaluation}}{\left(\text{Bias}(\hat{f}(\mathbf{x})) \right)^2} + \text{var}(\epsilon)$$

$$\text{Bias}(\hat{f}(\mathbf{x})) = E_T(\hat{f}(\mathbf{x}_t) - f(\mathbf{x}_t))$$

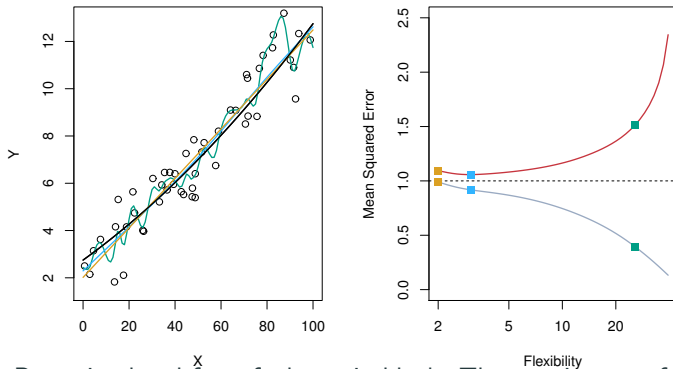
- [Lets derive the equation above]
- Normally as we make \hat{f} more flexible by making it complex (using more sophisticated formulations and more features) the **model variance** term $\text{var}(\hat{f}(\mathbf{x}_t))$ increases. This on the other hand decreases the bias (which we want to happen)
- So reducing the test error becomes a trade-off between the bias and the variance

Bias vs Variance



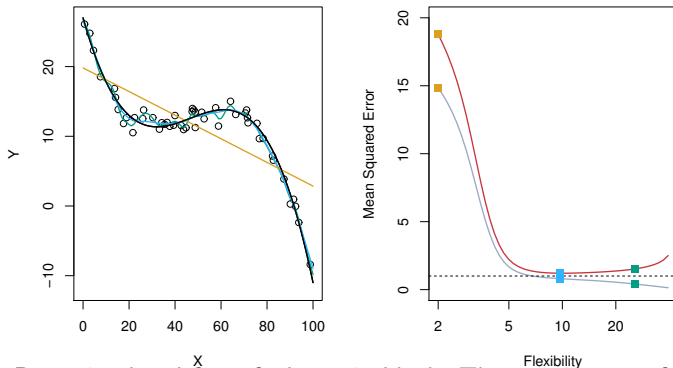
Left: Data simulated from f , shown in black. Three estimates of f are shown: the linear regression line (orange curve), and two smoothing spline fits (blue and green curves). Right: Training MSE (grey curve), test MSE (red curve), and minimum possible test MSE over all methods (dashed line). Squares represent the training and test MSEs for the three fits shown in the left-hand panel.

Bias vs Variance



Left: Data simulated from f , shown in black. Three estimates of f are shown: the linear regression line (orange curve), and two smoothing spline fits (blue and green curves). Right: Training MSE (grey curve), test MSE (red curve), and minimum possible test MSE over all methods (dashed line). Squares represent the training and test MSEs for the three fits shown in the left-hand panel.

Bias vs Variance



Left: Data simulated from f , shown in black. Three estimates of f are shown: the linear regression line (orange curve), and two smoothing spline fits (blue and green curves). Right: Training MSE (grey curve), test MSE (red curve), and minimum possible test MSE over all methods (dashed line). Squares represent the training and test MSEs for the three fits shown in the left-hand panel.

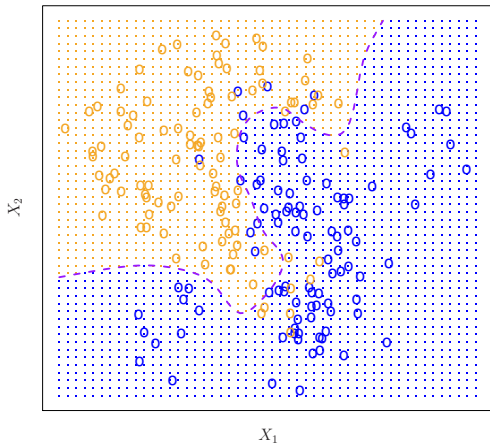
Trade Offs

We just discussed the bias vs variance trade off. In general we deal with various trade offs

- Prediction accuracy versus interpretability:
e.g., linear models are easy to interpret; thin-plate splines are not.
- Good fit versus over-fit or under-fit:
How do we know when the fit is just right? (just discussed on the figures above)
- Parsimony versus black-box:
We often prefer a simpler model involving fewer variables over a black-box predictor involving them all

Classification

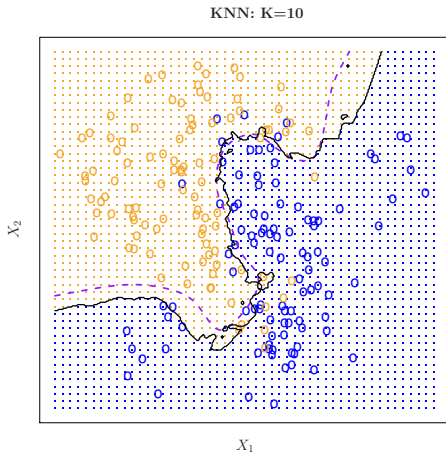
- How overfitting looks for classification problems
- K- Nearest Neighbors Approach



(ideal classifier)

Classification

- How overfitting looks for classification problems
- K- Nearest Neighbors Approach (the dense grid is somehow our test)

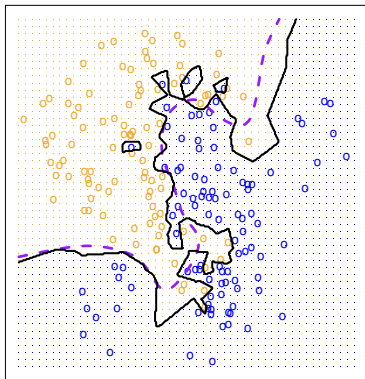


(KNN Classifier)

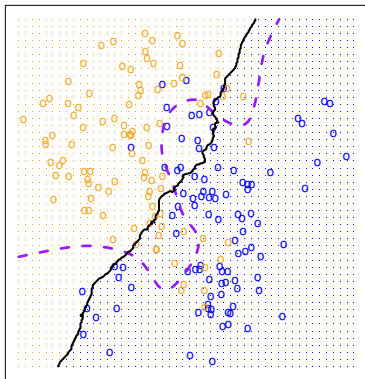
Classification

- How overfitting looks for classification problems
- K- Nearest Neighbors Approach (the dense grid is somehow our test)

KNN: K=1



KNN: K=100



(KNN Classifier)

Classification

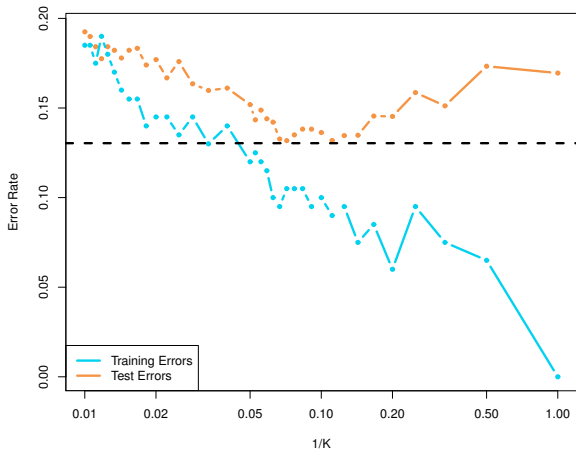
- K- Nearest Neighbors Approach
- Note that in KNN flexibility is inversely proportional to K



(Just like trying to walk among many people surrounding you)

Classification

- The Trade Off Curve
- K- Nearest Neighbors Approach



(Trade Off)

Lecture Overview

- Discussed the course information
- Related machine learning to model understanding in science and engineering
- Presented some machine learning examples in real world
- Talked about supervised vs unsupervised learning
- Talked about regression and classification
- Discussed some basics of statistical learning
 - Regression function
 - Test vs training data
 - Overfitting and model trade off
- **Don't worry if some of the topics covered look vague!** We will revisit all this material in depth later in the course

Bias ?