

ML in Practice

Mark Asch - ML-PREP

2025



Practical ML

Program

- Method choice: which method should I use for my problem?
 - ⇒ Bias-Variance tradeoff.
 - ⇒ Interpretability and explainability.
- EDA (Exploratory Data Analysis): how to use unsupervised learning as an initial phase?
- Model tuning and validation: how can I be sure of the accuracy and robustness of my model?
- Evaluation metrics and predictive power: which metrics measure the predictive power best?
 - ⇒ Dangers of R^2 and p -values.
- Causality vs Correlation: how to distinguish between them?

⇒ Shapley values

- Feature selection, feature importance and feature engineering.
- PINN (Physics Informed Neural Networks): when and how should they be incorporated?
- Ethics and bias.

CODING ENVIRONMENT

Setup

Python

Please follow general instructions in [01_setup.pdf](#)

ML FRAMEWORK

Mathematical Formulation

- Suppose that we have:

- ⇒ a **response** variable (to explain), Y ,
- ⇒ p **explanatory**¹ variables, $X = (X_1, X_2, \dots, X_p)$,
- ⇒ n **samples** of data, giving an $(n \times p)$ matrix,

$$X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & & \\ \vdots & & \ddots & \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix}$$

- ⇒ a relationship between Y and X of the form

$$Y = f(X) + \epsilon$$

- ⇒ where

→ f is an **unknown** function of (columns)
 X_1, X_2, \dots, X_p

¹Also called: **features**, **attributes**, **covariates**

→ ϵ is a random **error** term, independent of X , and with zero mean

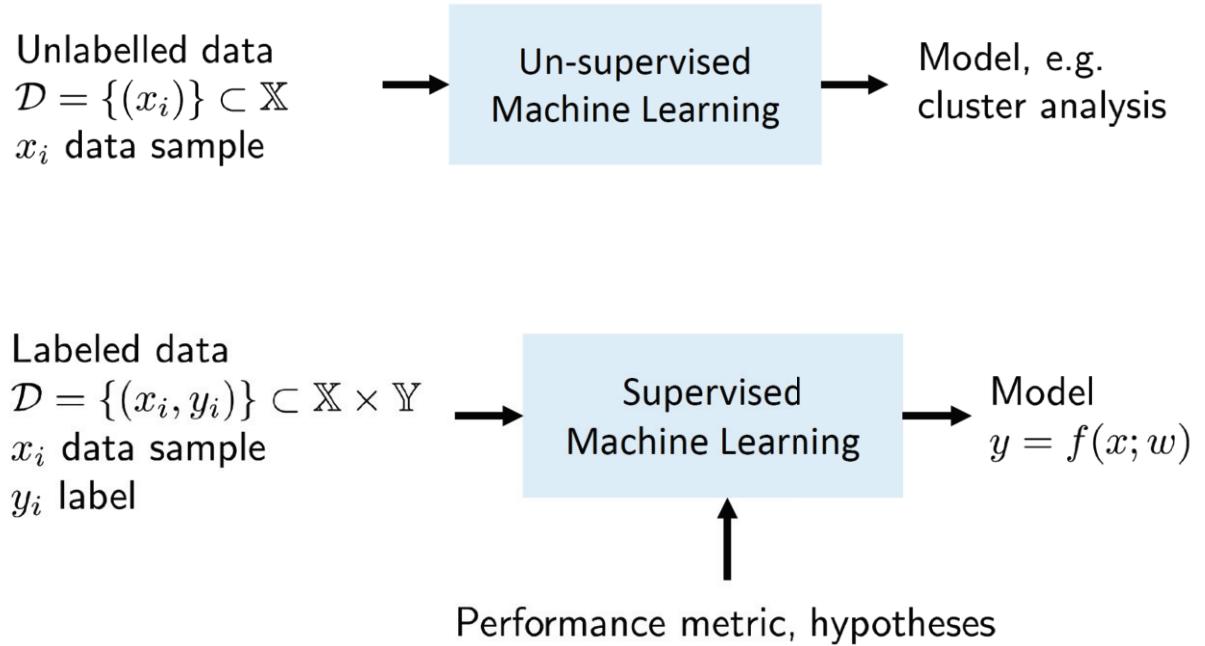
- ML is then an ensemble of approaches for **estimating** the unknown function f with the objectives of
 - ⇒ **Prediction**: $\hat{Y} = \hat{f}(X)$ where \hat{f} is an estimation for f and \hat{Y} is the resulting prediction
 - ⇒ **Inference**: to understand how Y varies as a function of X (correlations, importances, linearity, etc.)

Initial Categories

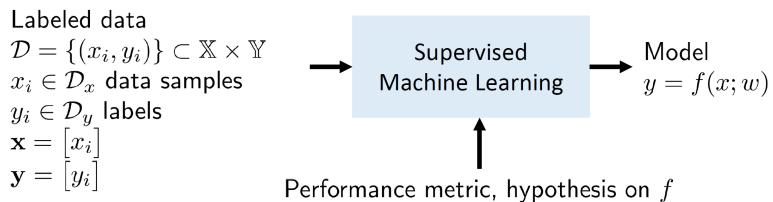
See: [Overview - general](#), [Overview - methods](#)

1. Supervised vs. Unsupervised.
2. Regression vs. Classification.
3. Mixtures of the above?

Supervised vs. Unsupervised

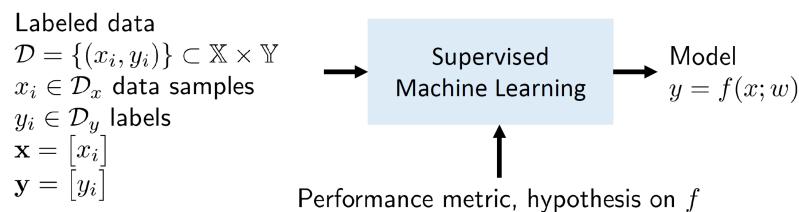


Regression vs. Classification



Model purpose – Regression

- The model f shall map $x \mapsto y$ and approximate an unknown function
 $\hat{f} : \mathbb{X} \rightarrow \mathbb{Y}$
- $y_i \in \mathbb{Y} \subseteq \mathbb{R}^{n_y}$
- Examples: data-driven modeling, energy forecasting, ...



Model purpose – Classification

- The model f shall map $x \mapsto y$ and approximate an unknown function
 $\hat{f} : \mathbb{X} \rightarrow \mathbb{Y}$
- $y_i \in \mathbb{Y} \subseteq \mathbb{N}^{n_y}$
- Examples: spam filter, fraud detection, fault detection, ...

The only difference is the space over which the response Y is defined!

METHOD CHOICE

Classes of ML Methods

- supervised learning
- unsupervised learning
- reinforcement learning (AlphaGo, AlphaFold)
- self-supervised learning (LLMs, ChatGPT)

Table of ML Methods

Class	Model	Task
Supervised	Linear regression	R
	CART (trees)	R, C
	SVM	R, C
	NN	R, C
	k -NN	C
	Naive Bayes	C
Unsupervised	k -means	clustering
	hierarchical	clustering
	PCA	patterns

- R = regression, C = classification
- CART = Classification and Regression Trees
- SVM = Support Vector Machine
- NN = Neural Network
- PCA = Principal Component Analysis

Methods - presentations

Lectures and code examples for all the above methods are available on the CMU-IMU-2023 [course website](#).

How to choose a model?

- Once the **tuning** (hyper)parameters have been determined/chosen (see: cross-validation below), we still must choose between several models
 - ⇒ the choice will largely depend on the data characteristics and the type of questions we ask
- But, **predicting** which model will be the most pertinent, is in general quite difficult...

Scheme for “How to choose a model”

A recommended scheme for finalizing the choice is as follows:

1. Begin with a few models that are the least interpretable and the most flexible. These models have a strong chance of producing better precision.
 - (a) SVM.
 - (b) Random Forests (RF).
 - (c) Trees with boosting (xxxBoost).
2. Study simpler models, that are less opaque.
 - (a) Linear models.
 - (b) (GAM/GLM)
 - (c) Naive Bayes.
 - (d) k-NN.
 - (e) Logistic Regression.
 - (f) Regression Splines (MARS).

3. Use, if possible, the **simplest** model (see: bias-variance trade-off) that approximates reasonably well the performance of the more **complex** models.

Remark 1. In ML, the algorithm learns the **parameters** of the model, usually by optimization of a error-loss score. The method itself is defined by **hyperparameters**. These hyperparameters control the learning process. They have default values, but need to be **tuned** for each new dataset.

EDA

Recall: EDA

See: [Overview - general](#)

Remark 2. In addition to plots and summary statistics, [unsupervised learning](#) techniques often play an important role in identifying structure in the data.

CROSS-VALIDATION and TUNING

Recall: CV and Tuning

See: [Resample](#), [elementary train-test tutorial](#), a bit more advanced [CV tutorial](#)

- Take-home lessons:
 - ⇒ Basic, [train-test split](#) is not recommended since a single split may not be representative, and the resulting score may not be a reliable estimate of the predictive power on unseen data (i.e. in production.)
 - ⇒ Simple, [k-fold CV](#) is better, but is rarely sufficient to ensure model reliability.
 - ⇒ [Nested CV](#) with inner loop to select hyperparameters, outer loop to evaluate the best model, is probably the best option.
 - ⇒ [Repeated k-fold CV](#).
 - A single run of the k-fold cross-validation procedure may result in a noisy estimate of model performance.
 - Different splits of the data may result in very different results.

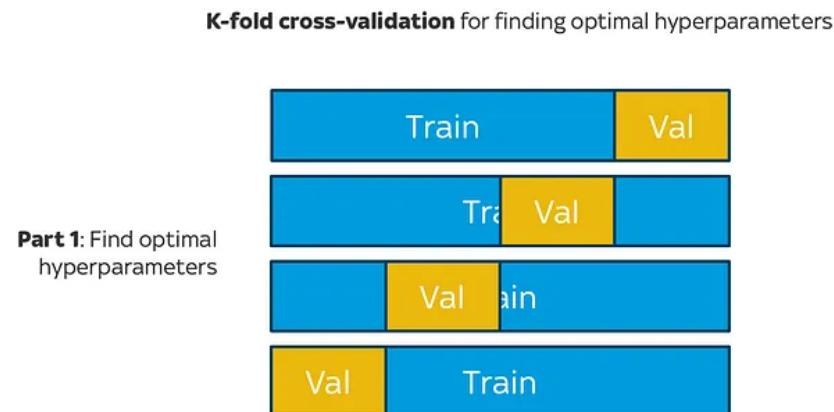
- The **mean result** is expected to be a more accurate estimate of the true unknown underlying mean performance of the model on the dataset, as calculated using the standard error.
- Usual numbers of repeats include 3, 5, and 10.
- ⇒ **Time-series data** requires special splitting, that respects the temporality (autocorrelation). For this one should use the scikit function **TimeSeriesSplit**
- ⇒ **Geospatial data** often requires special splitting, that respects the spatial distribution (autocorrelation). For this one should use special techniques, presented in the advanced geospatial data analysis lecture notes.

Nested CV

Single split:



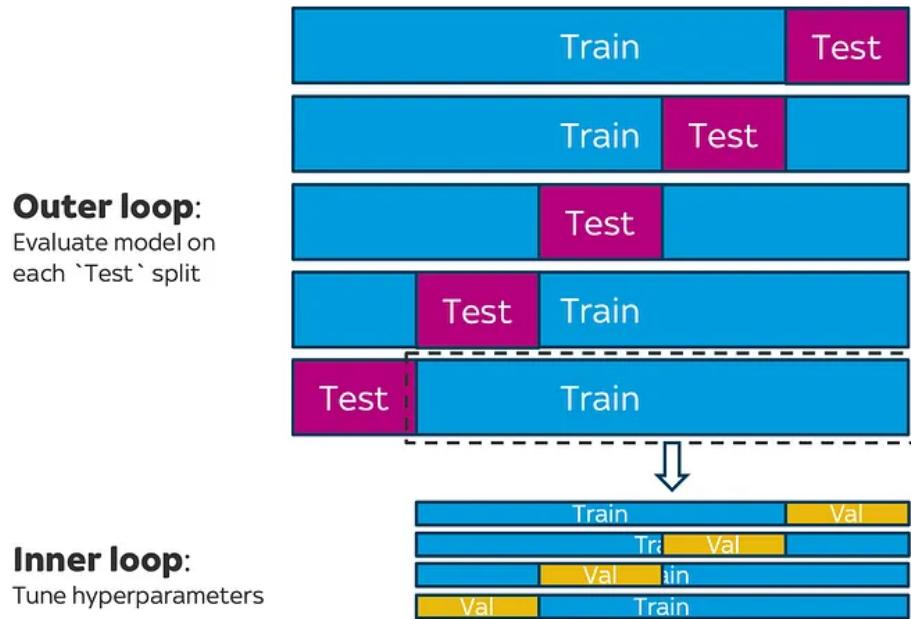
k -fold CV:



Part 2: Using optimal hyperparameters found in Part 1, train a model on the entire training split, and evaluate on the held-out test set*



Nested k -fold CV:



See Example notebook.

Resampling: Take-Home Lessons

No resampling method is uniformly better than another and the choice of a particular method should consider several aspects.

- If the **sample size is small**, repeated (5- or) 10-fold cross-validation is recommended for several reasons:
 - ⇒ the bias and variance properties are good and,
 - ⇒ given the sample size, the computational costs are not large.
- If the goal is to **choose between models**, as opposed to getting the best indicator of performance, it is advisable to use a bootstrap method since these have very low variance.
- For **large sample sizes**, the differences between resampling methods become less pronounced, and computational efficiency becomes the critical factor. Here,

simple 10-fold cross-validation should provide acceptable variance and low bias, and is relatively cheap to compute.

Finalized Model

- A **final** machine learning model is a model that you will use to make predictions on new data.
- You finalize a model by applying the chosen, tuned machine learning procedure on **all of your data**. That's it.
- See Example notebook.

EVALUATION METRICS

Loss Functions and Metrics

We need to distinguish between metrics for the 2 following cases. See: [linear regression](#), [model selection](#).

- For a [given model](#): MSE, R^2 , etc.
- For comparing [different models](#): AIC, BIC, etc.

Definitions

- Coefficient of determination (R-squared)

$$R^2 = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2} = 1 - \frac{\text{RSS}}{\text{TSS}} = \frac{\text{TSS} - \text{RSS}}{\text{TSS}}$$

represents the proportion of the variance “explained” by the model. Beware of the **limitations** of this coefficient!

- Mean-squared Error:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2,$$

where \hat{y}_i is the predicted value by the model for observation i . This can be more reliable than R-squared, but is context-dependent (not between 0 and 1).

- Mean Classification Error:

$$\text{MCE} = \frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i),$$

where I is the indicator function and \hat{y}_i is the class, or label predicted by the model for observation i .

Recall: Measures of Performance

- We compute **individual** performance scores for each cycle of the CV loop:
 - ⇒ For regression problems, use the **mean squared error**.
 - ⇒ For classification problems, use the **misclassification rate**.
- Then, the **general score** of the CV is the average of the k individual scores,

$$\text{CV}_{(k)} = \frac{1}{k} \sum_{j=1}^k \text{MSE}_j,$$

to which we add a measure of the variance, the **standard error**, which is the standard deviation of the sampling distribution,

$$\text{SE} = \frac{s}{\sqrt{n}},$$

where s is the empirical standard deviation, the square root of the empirical variance,

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2.$$

- Finally, an estimated 95% **confidence interval** of the empirical mean², is defined by

$$\bar{x} \pm 1.96\text{SE}.$$

²This confidence interval is exact for a Gaussian distribution, but random i.i.d. errors and their sums are only approximately Gaussian.

Recall: Mean-Squared Error

- Since we are estimating an unknown parameter, or more often, unknown function, we need to evaluate
 - ⇒ how good our estimation is?
 - ⇒ the quality of the estimation?
- This is traditionally done by using a discrete L_2 norm, though other norms are possible and are used in specific circumstances—see below.
- Before defining the MSE, we need to recall the definitions of
 - ⇒ expectation (or expected value):

$$\mathbb{E}[X] = \sum_i x_i f(x_i),$$

where f is the (discrete) probability function.

⇒ bias:

$$\text{Bias}(T, \theta) = \text{Bias}(T) = E[T] - \theta,$$

where T is a statistic used to estimate a parameter, and θ is the parameter to be estimated

⇒ variance:

$$\text{Var}[X] = E[(X - \mu)^2] = E[X^2] - \mu^2,$$

where μ is the mean of X .

Definition 1 (Mean Squared Error). The mean squared error measures the **average of the squares of the errors**—that is, the average squared difference between the estimated values and the actual value. MSE is a risk function, corresponding to the expected value of the squared error loss,

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{1}{n} \sum_{i=1}^n (e_i)^2,$$

or for an estimator

$$\text{MSE}(\hat{\theta}) = \text{E}_{\theta} \left[(\hat{\theta} - \theta)^2 \right].$$

The MSE is the second moment (about the origin) of the error, and thus incorporates both the **variance** of the estimator (how widely spread the estimates are from one data sample to another) and its **bias** (how far off the average estimated value is from the true value).

The MSE either assesses the quality of a **predictor** (i.e., a function mapping arbitrary inputs to a sample of values of some random variable), or of an **estimator** (i.e., a mathematical function mapping a sample of data to an estimate of a parameter of the population from which the data is sampled). The definition of an MSE differs according to whether one is describing a predictor or an estimator.

Remark 3. The mean squared error has the disadvantage of heavily weighting **outliers**. This is a result of the squaring of each term, which effectively weights large errors more heavily than small ones. This property, undesirable in many

applications, has led to the use alternatives such as the **mean absolute error**, which is a discrete L_1 norm, or scores based on the **median**.

Definition 2 (Mean Absolute Error). The mean absolute error is the average of the sum of absolute errors,

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| = \frac{1}{n} \sum_{i=1}^n |e_i| ,$$

or for an estimator

- A python snippet:

```
from sklearn.metrics import mean_absolute_error
import numpy as np
# Generate some sample data
y_true = np.array([1, 2, 3, 4, 5])
y_pred = np.array([1.5, 2.5, 2.8, 4.2, 4.9])
# Calculate the MAE
mae = mean_absolute_error(y_true, y_pred)
print("Mean Absolute Error", mae)
```

Remark. The MAE has two important properties.

1. It is **robust** to outliers. The MAD (**median** absolute deviation) is even more robust.
2. It has the same **units** as the response variable.

Bias-Variance Decomposition of MSE

Theorem 1. *The bias-variance decomposition of the MSE is given by*

$$\begin{aligned}\text{MSE} &\doteq \mathbb{E} \left[(y - \hat{f})^2 \right] \\ &= \text{Bias} \left[\hat{f} \right]^2 + \text{Var} \left[\hat{f} \right] + \sigma^2,\end{aligned}$$

where $y = f(x) + \epsilon$, $\epsilon \sim \mathcal{N}(0, \sigma^2)$ ³ and $\hat{f}(x)$ is an approximation/model of the true (unknown) function $f(x)$.

Proof. First expand the definition of MSE,

$$\begin{aligned}\text{MSE} &\doteq \mathbb{E} \left[(y - \hat{f})^2 \right] \\ &= \mathbb{E} \left[y^2 - 2y\hat{f} + \hat{f}^2 \right] \\ &= \mathbb{E} [y^2] - 2\mathbb{E} [y\hat{f}] + \mathbb{E} [\hat{f}^2].\end{aligned}$$

³This Gaussian noise model is the most arbitrary (least restrictive), assuming only a central tendency and low chance of extreme values.

Now look at each of the three terms.

$$E[\hat{f}^2] = \text{Var}[\hat{f}] + E[\hat{f}]^2$$

from variance decomposition.

$$\begin{aligned} E[y^2] &= E[(f + \epsilon)^2] \\ &= E[f^2] + 2E[f\epsilon] + E[\epsilon^2] \\ &= f^2 + 2f \cdot 0 + \sigma^2 \end{aligned}$$

from definition of the noise distribution ϵ . Finally,

$$\begin{aligned} E[y\hat{f}] &= E[(f + \epsilon)\hat{f}] \\ &= E[f\hat{f}] + E[\epsilon]E[\hat{f}] \\ &= fE[\hat{f}]. \end{aligned}$$

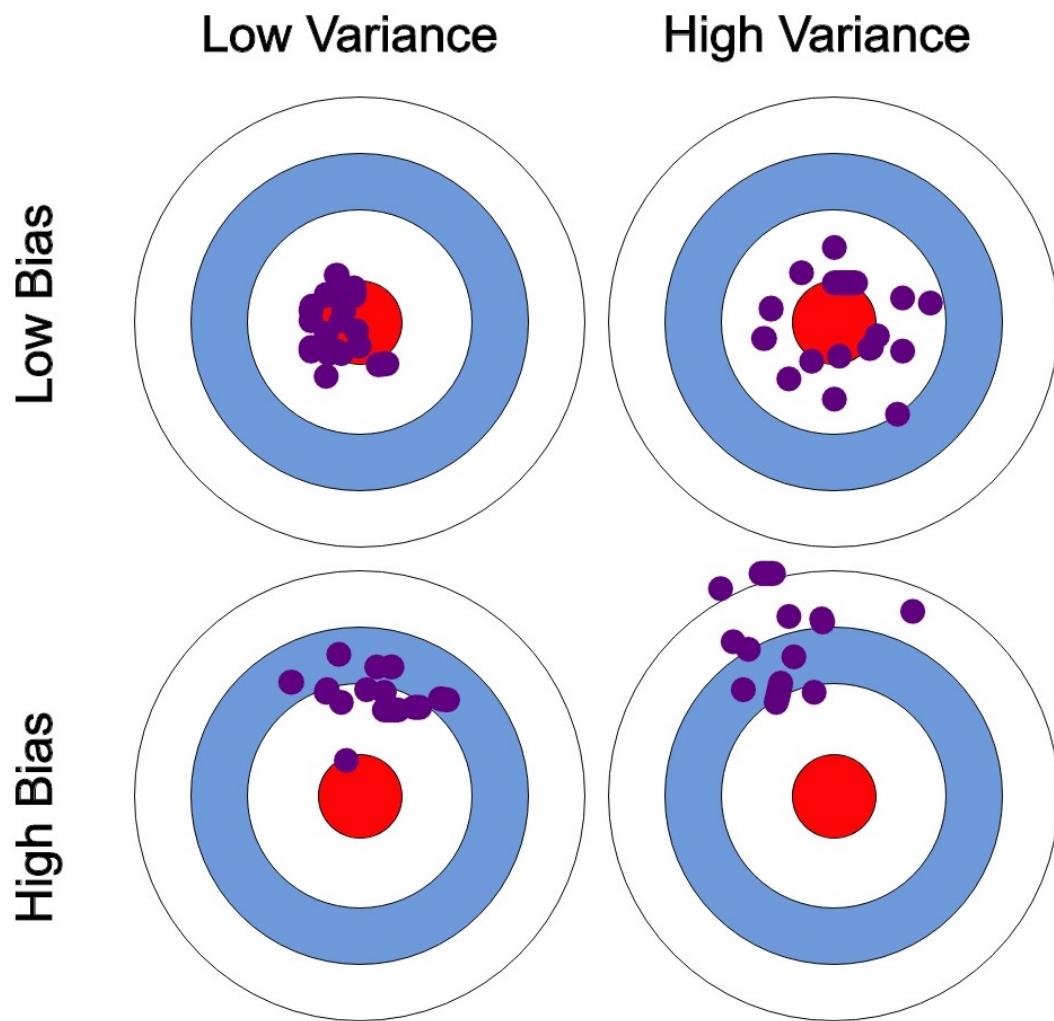
Putting it all together, we find

$$\begin{aligned}\text{MSE} &= (f^2 + \sigma^2) - 2f\mathbb{E}[\hat{f}] + \left(\text{Var}[\hat{f}] + \mathbb{E}[\hat{f}]^2\right) \\ &= \left(f - \mathbb{E}[\hat{f}]\right)^2 + \sigma^2 + \text{Var}[\hat{f}] \\ &= \text{Bias}[\hat{f}]^2 + \text{Var}[\hat{f}] + \sigma^2.\end{aligned}$$

□

Recall: Bias-Variance Compromise

- Questions:
 - ⇒ What is bias? (see also **Ethics** lecture)
 - **Low Bias**: Predicted data points are close to the target. The model has fewer assumptions about the form of the target function.
 - **High-Bias**: Predicted data points are far from the target. The model has more assumptions about the form of the target function.
 - ⇒ What is variance?
 - **Low Variance**: Data points are close to each other, and as a result are close to the function. The model suggests small changes to the estimate of the target function with changes in the training dataset.
 - **High Variance**: Data points are spread out and as a result are far from the function. The model suggests large changes to the estimate of the target function with changes in the training dataset.



Bias-Variance Compromise II

- We can now return to the question of the **tradeoff** between bias and variance.
- Looking at the bias-variance decomposition of the MSE,

$$\text{MSE} = \text{Bias} [\hat{f}]^2 + \text{Var} [\hat{f}] + \sigma^2,$$

we observe that we need to **simultaneously** minimize both bias and variance. This is the **reducible** part of the error, whereas σ^2 is the intrinsic, **irreducible** part and cannot be reduced by the ML model.

- Theoretical result (see also the “compromise curve” above):
 - ⇒ Decreasing the bias will mechanically increase the variance.
 - Low bias \implies high variance.

- ⇒ Decreasing the variance will mechanically increase the bias.
 - Low variance \implies high bias.

- **Conclusion:**

- ⇒ To build a good model, we need to find a good **balance/compromise/tradeoff** between bias and variance that minimizes the total/generalization error.
- ⇒ An optimal balance of bias and variance should neither overfit nor underfit the model. This is achieved by **hyperparameter tuning** on the basis of systematic grid-search (trial and error, or experience, for the ranges) and cross-validation. See above, and examples.

A simple example that brings out many points

- **Basic fact:** most of ML is about **pattern recognition**—see C. M. Bishop, PRML, Springer, 2006.
 - ⇒ we have a bunch of data, \mathbf{x} and corresponding observations, \mathbf{y}
 - ⇒ we seek the relation (recognize the pattern), f , that relates \mathbf{x} to \mathbf{y}

$$\mathbf{y} = f(\mathbf{x}),$$

or,

$$f: \mathbf{x} \longmapsto \mathbf{y}$$

- **Our goal:** make **predictions** for new (unseen) data (inputs)

Polynomial Curve Fitting

- suppose that

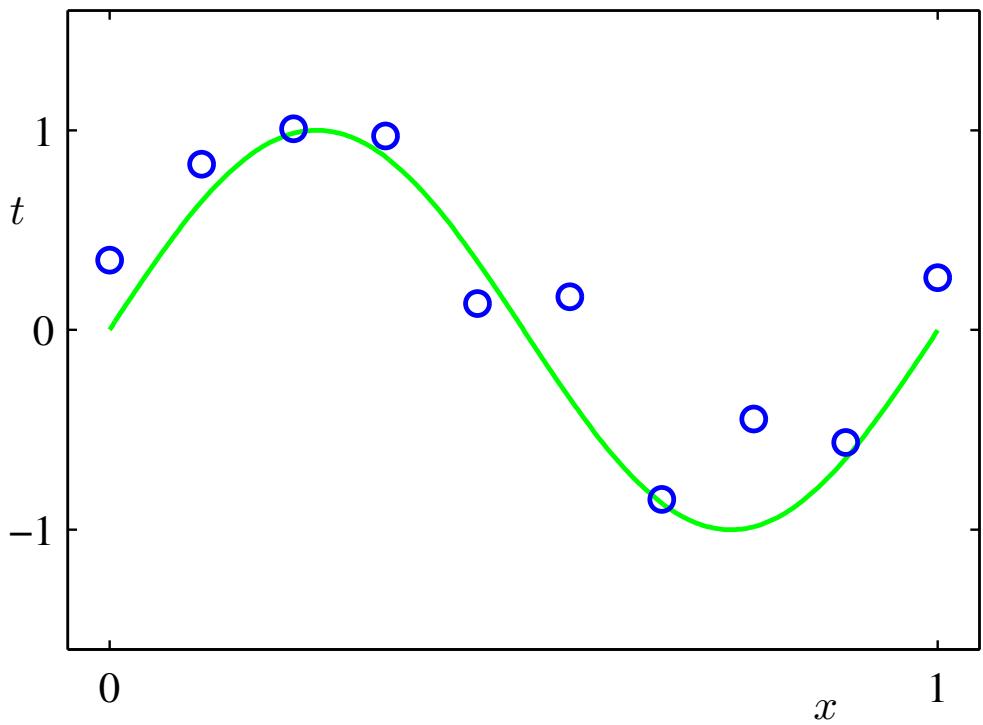
$$y = \sin(2\pi x) + \text{noise}$$

- given: a **training** set of N observations of x ,

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$$

and the corresponding (noisy) **observations**,

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$



- data generation from

$$y_i = \sin(2\pi x_i) + \epsilon,$$

where the noise term,

$$\epsilon \sim \mathcal{N}(0, \sigma)$$

is due to intrinsic/natural variability and other uncertainties from the data collection methods.

- **Objective:** from the training set, learn/discover an approximation of the relation f that will enable us to make predictions based on
 - ⇒ probability theory
 - ⇒ (polynomial) curve fitting
- PCF: suppose that the data can be represented by the **polynomial** / "linear model"

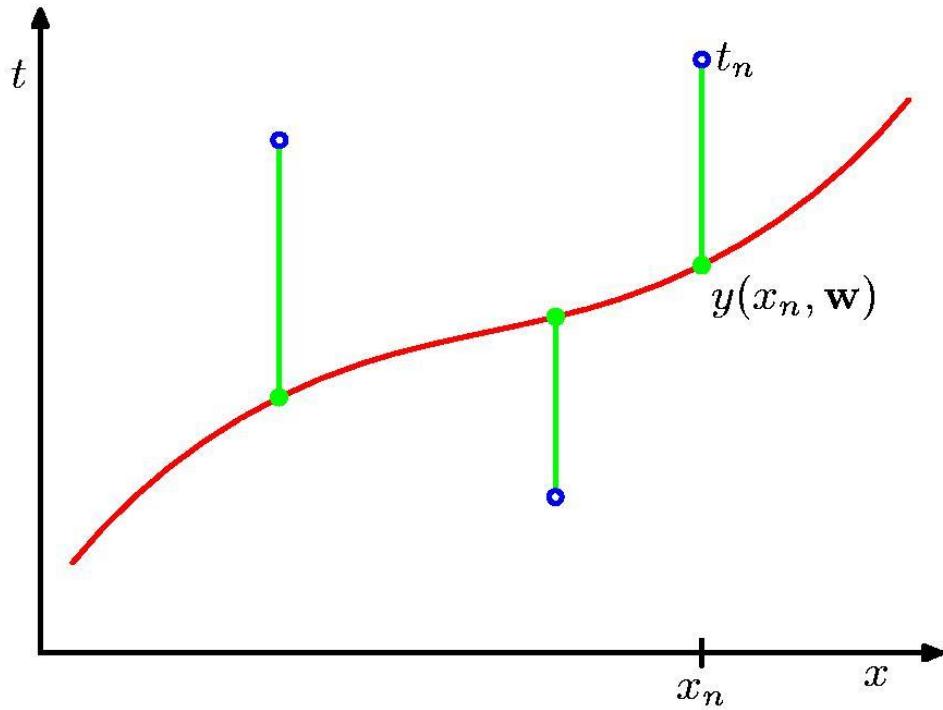
$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \cdots + w_M x^M = \sum_{k=0}^M w_k x^k$$

with **weight** vector (coefficients) \mathbf{w} .

⇒ compute the coefficients by minimizing an **error/loss function**—this is (often) a least-squares approach, where the loss function is defined as

$$\mathcal{L}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N [\hat{y}(x_i, \mathbf{w}) - y_i],$$

where \hat{y}_i is the approximation obtained from our polynomial model.



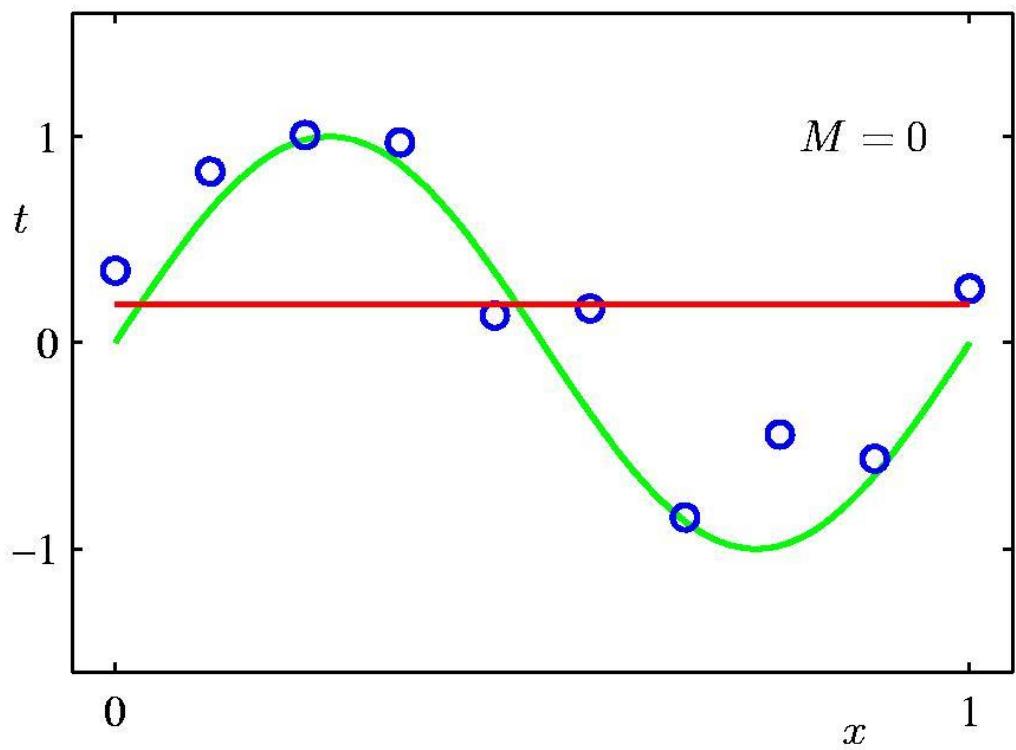
- in the case of a quadratic loss function, the gradient is linear and the solution to the optimization problem exists and can be found in closed form using the normal form, or psuedo-inverse,

$$\mathbf{w}_* = (X^T X)^{-1} X^T \mathbf{y},$$

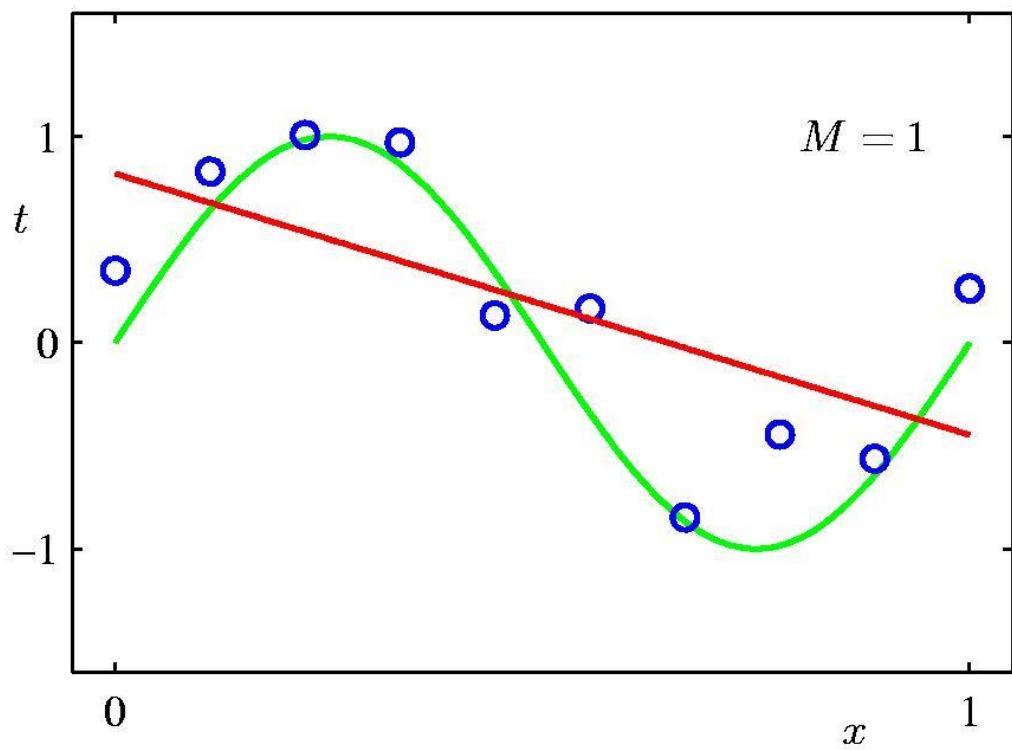
where X is the (rectangular) coefficient matrix of dimension $(N \times M)$ with $N > M$ and \mathbf{y} is the data/observation vector.

- the choice of M is a **model selection** problem
 - ⇒ we compute the solution for 4 values: $M = 0, 1, 3, 9$ —see below
 - ⇒ Model $M = 0$ produces **underfitting**—it is too simple and does not find the trend
 - ⇒ Model $M = P$ produces **overfitting**—it fits the data points perfectly, but is “brittle” and does not generalize
- **Conclusion:** we must seek a **compromise** between the two.

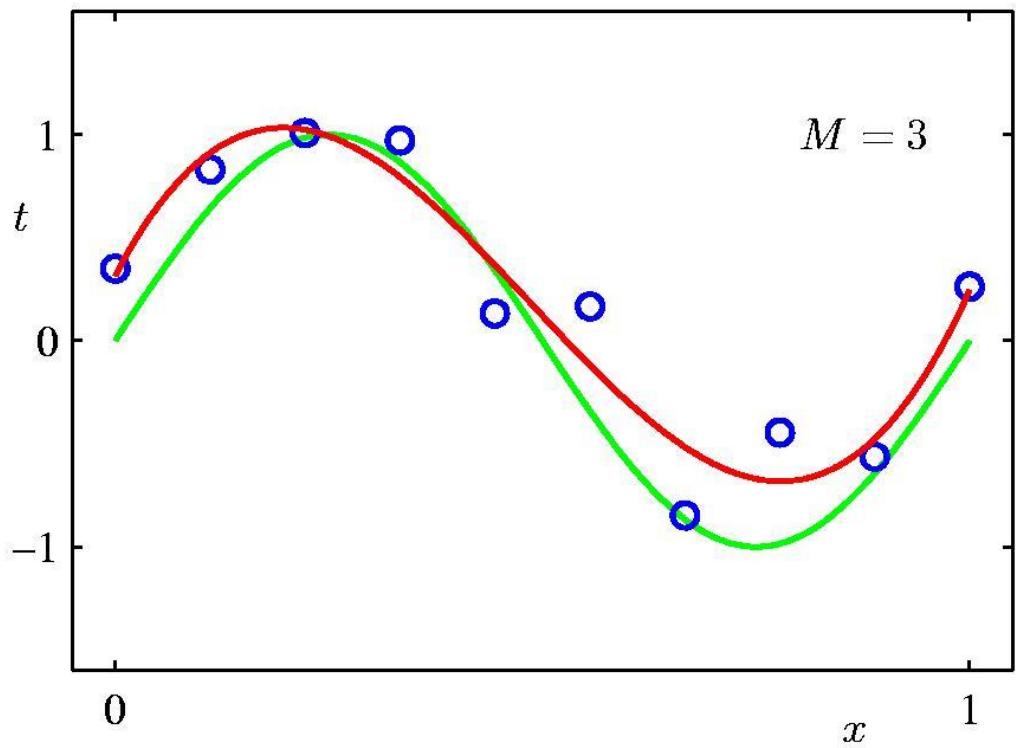
0th Order Polynomial



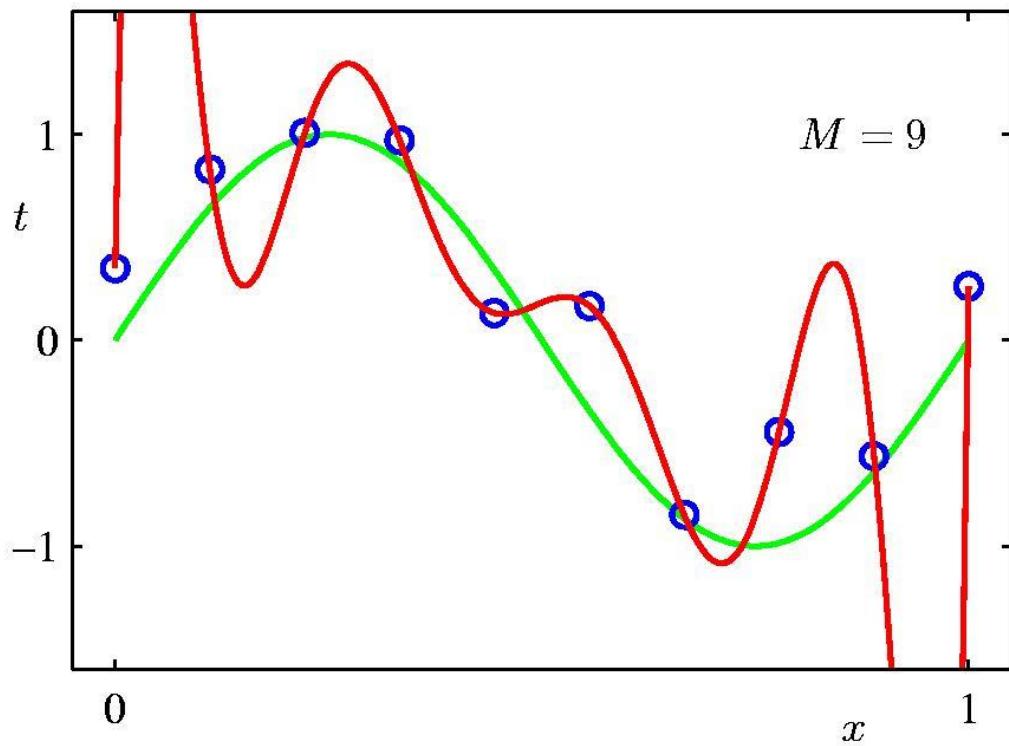
1st Order Polynomial



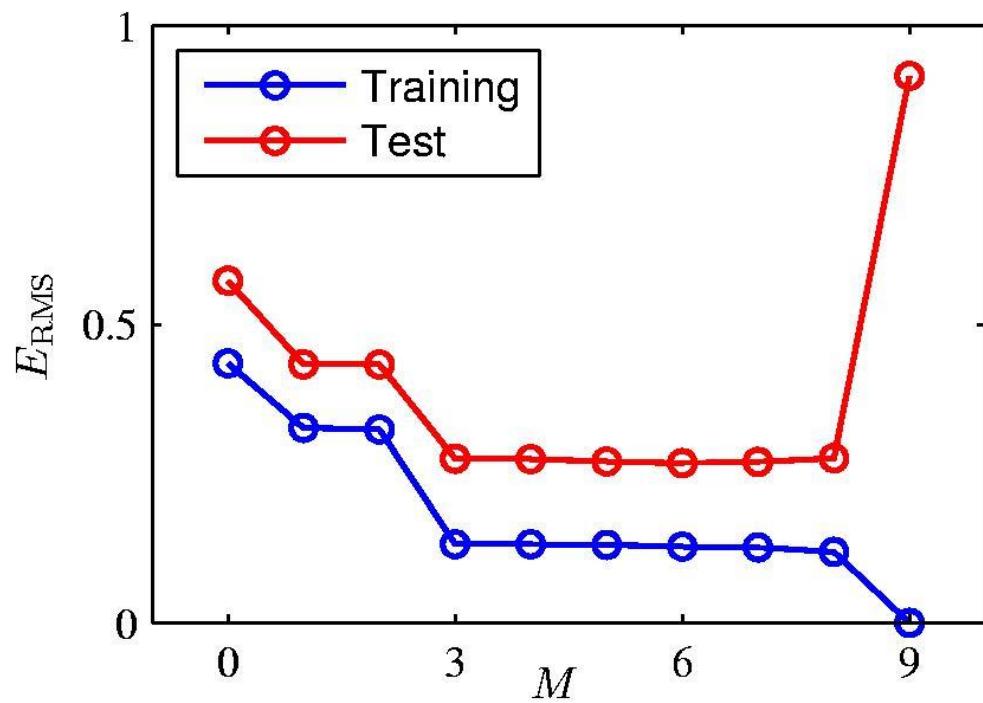
3rd Order Polynomial



9^{th} Order Polynomial



Over-fitting

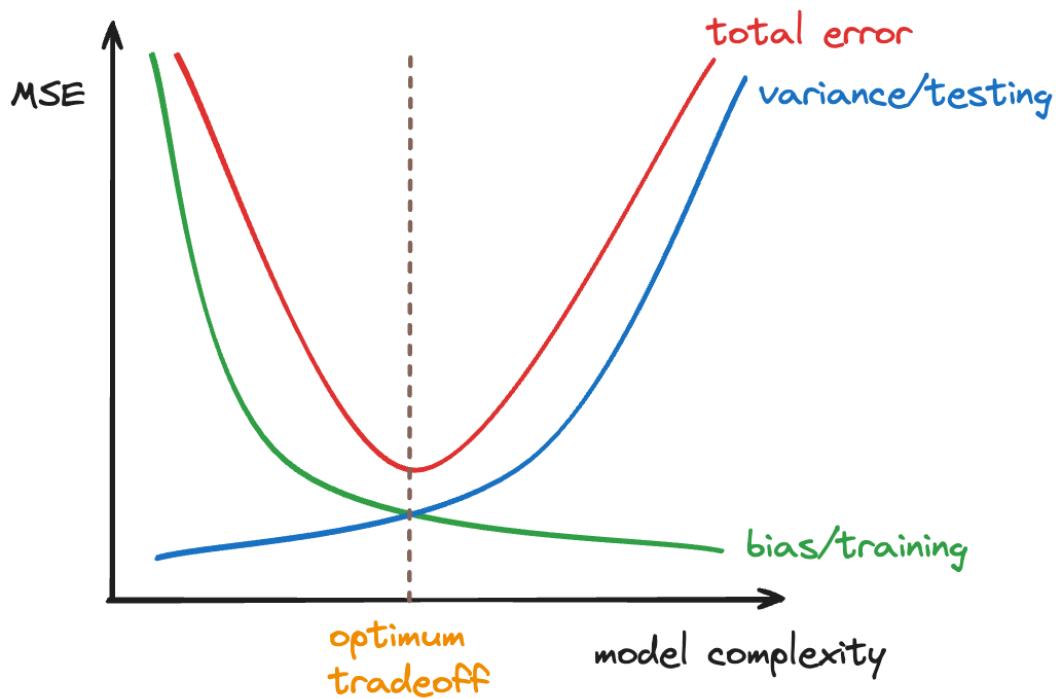


Root-Mean-Square (RMS) Error: $E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$

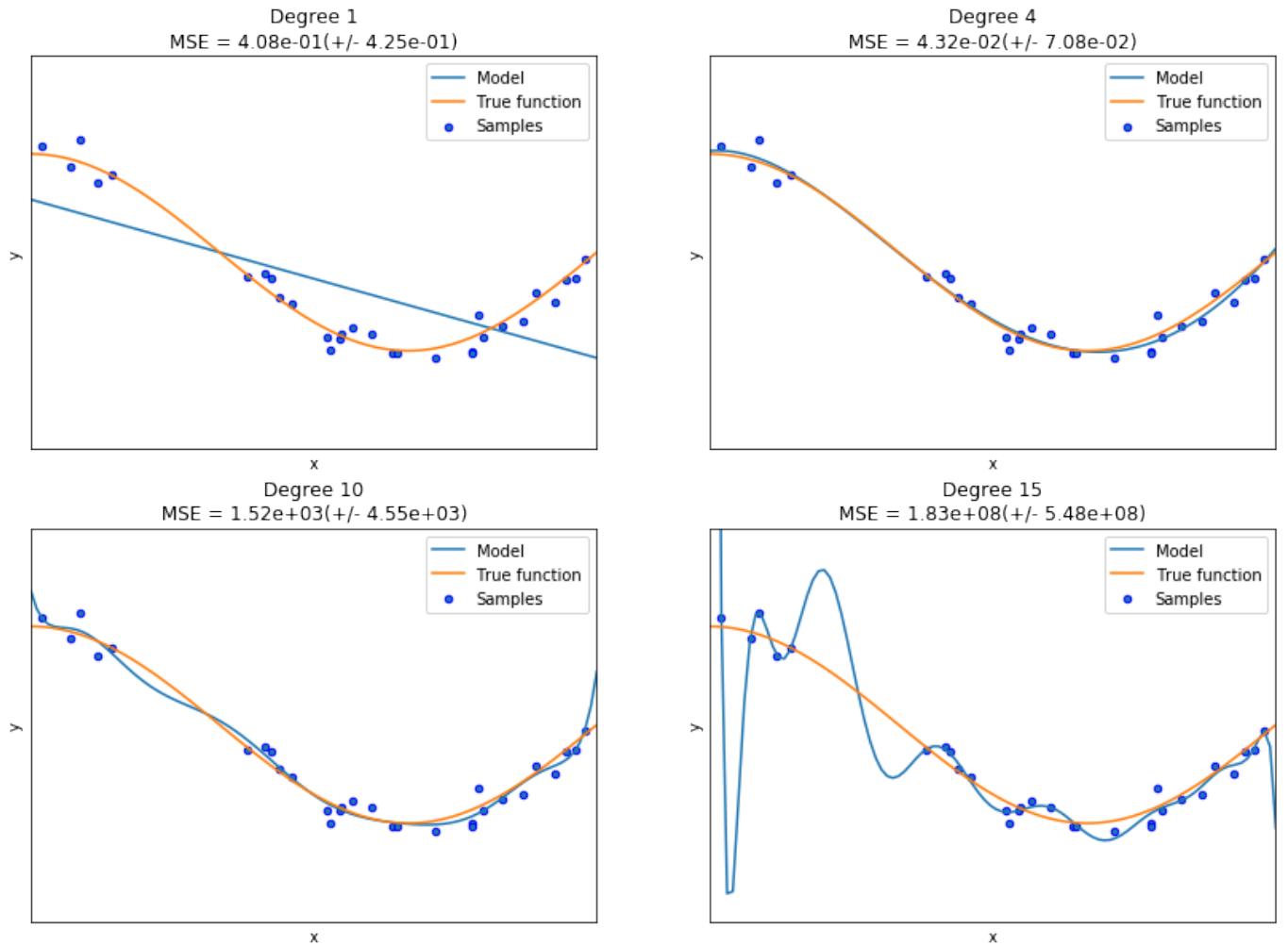
Polynomial Coefficients

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43

Bias-Variance Compromise



- We saw the dangers of **overfitting** in the very simple polynomial fitting example.
- To treat this, we need to make a **compromise** between bias and variance, i.e. forgo some precision in order to attain better **predictive performance**.
⇒ this is a fundamental point underlying **ALL** machine learning methods



order M	MSE
1	0.408
4	0.043
10	$1.52e03$
15	$1.83e08$

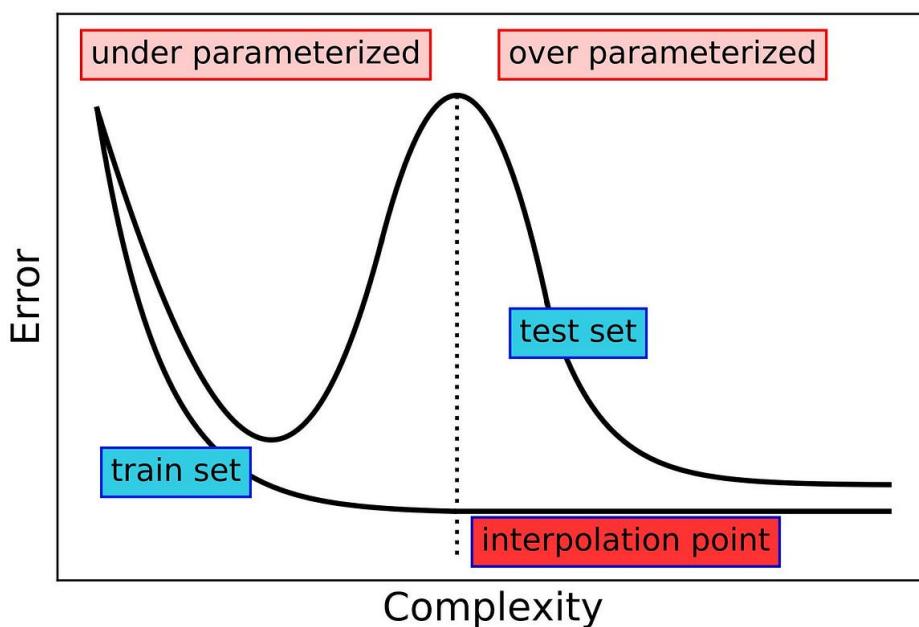
Over- and Underfitting

- ML Algorithm properties:
 - ⇒ Examples of **low-bias** machine learning algorithms: Decision Trees, k-Nearest Neighbors and Support Vector Machines.
 - ⇒ Examples of **high-bias** machine learning algorithms: Linear Regression, Linear Discriminant Analysis and Logistic Regression.
 - ⇒ Examples of **low-variance** machine learning algorithms: Linear Regression, Linear Discriminant Analysis and Logistic Regression.
 - ⇒ Examples of **high-variance** machine learning algorithms: Decision Trees, k-Nearest Neighbors and Support Vector Machines.
- Dealing with the overfitting problem:
 - ⇒ always be prepared to “sacrifice” some bias, to get a lower variance, since it is the variance that is the most “dangerous” in terms of **predictive power/performance**.

- ⇒ in general, prefer lower-order ML methods—so-called **linear methods**
- ⇒ increase the number of **data points**—this is not always feasible in reality
- ⇒ use a **Bayesian** approach—probably the best, but more complex to implement
- ⇒ **regularize** the problem—penalize terms of large magnitude, but requires extensive tuning of the regularization parameters that usually relies on **cross-validation** techniques.

BV Tradeoff: Double Descent phenomenon

- Test error can present a “double descent” phenomenon in a range of ML models.
- As the number of model parameters grows relative to the number of data points, test error drops in the highly **overparameterized** (data undersampled) regime



- Interpretation: [OpenAI]
 - ⇒ For models at the interpolation threshold, there is effectively only one model that fits the train data, and forcing it to fit even slightly noisy or misspecified labels will destroy its global structure.
 - ⇒ There are no “good models” which both interpolate the train set and perform well on the test set.
 - ⇒ However, in the over-parameterized regime, there are many models that fit the train set and there exist such good models.
 - ⇒ Moreover, the implicit bias of stochastic gradient descent (SGD) leads it to such good models, for reasons we do not yet understand...
 - ⇒ **Warning:** this is a topic of ongoing research!

Recall: Precision-Recall Compromise

- What happens with **binary classification** problems?
 - ⇒ the famous “truth/confusion table”:
 - True positives (TP): positive COVID test, correctly diagnosed.
 - False positives (FP): positive COVID test, falsely diagnosed.
 - True negatives (TN): negative COVID test, correctly diagnosed.
 - False negatives (FN): negative COVID test, falsely diagnosed.

		True Class	
		Positive	Negative
Predicted Class	Positive	TP	FP
	Negative	FN	TN

Definition 3. A *confusion table/matrix*, C , for a classification with n classes is an $(n \times n)$ matrix with entries

C_{ij} = the number of observations in class i
that are predicted in class j .

Then,

- C_{ii} , $i = 1, \dots, n$, are the *good* classifications,
- C_{ij} with $i \neq j$ are the *bad* classifications.
- Now we can define precision and recall/sensitivity.

Precision the number of true positives divided by the total number of true positives and false positives:

$$\text{PPV} = \frac{\text{TP}}{\text{TP} + \text{FP}} = 1 - \text{FDR}$$

Appropriate when minimizing false positives is the focus—FP is more serious, think of a false positive cancer diagnosis, and the heavy, unnecessary treatment that follows.

Recall the number of correct positive predictions made out of all positive predictions that could have been made:

$$\text{TPR} = \frac{\text{TP}}{\text{P}} = \frac{\text{TP}}{\text{TP} + \text{FN}} = 1 - \text{FNR}$$

Appropriate when minimizing false negatives is the focus—FN is more serious, think of missing a positive COVID, who then goes on to infect many people.

Specificity is the number of true negatives divided by the total number of true negatives and false positives:

$$\text{TNR} = \frac{\text{TN}}{\text{N}} = \frac{\text{TN}}{\text{TN} + \text{FP}} = 1 - \text{FPR}$$

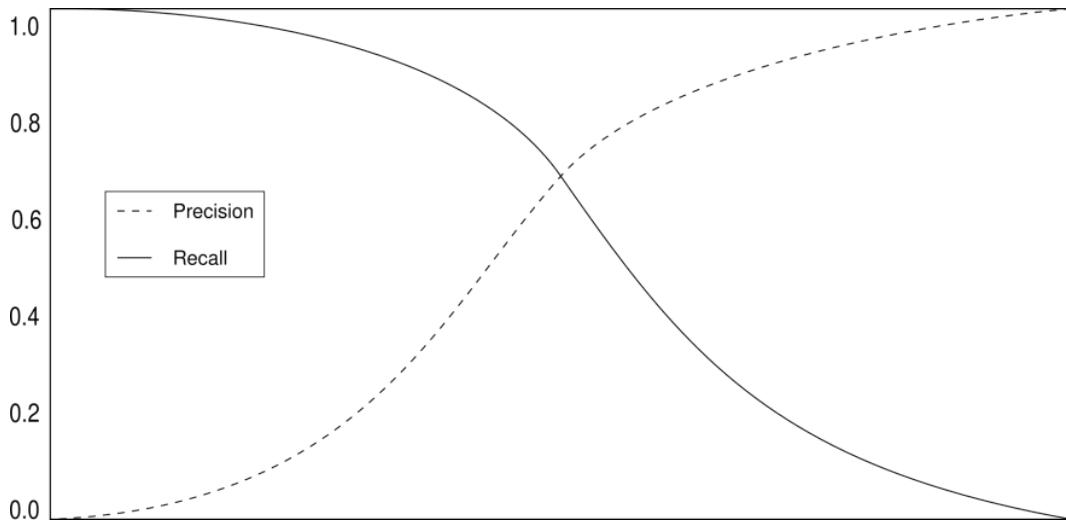
F1-Score is the harmonic mean of precision and recall:

$$F_1 = 2 \times \frac{\text{PPV} \times \text{TPR}}{\text{PPV} + \text{TPR}} = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}}$$

Alone, neither precision or recall tells the whole story. We can have excellent precision with terrible

recall, or alternately, terrible precision with excellent recall. F-measure provides a way to express both concerns with a single score.

Precision-Recall Compromise II



- The **compromise**:
 - ⇒ When precision increases, recall decreases, and vice-versa.
 - ⇒ We cannot simultaneously increase both precision and recall.
- To balance precision and recall, a number of **techniques** can be used, such as
 - ⇒ adjusting the decision threshold or

- ⇒ using an ensemble of models.
- Another approach is to use a metric that combines both precision and recall,
 - ⇒ such as the F1 score, which is the harmonic mean of precision and recall.
- Ultimately, the choice between precision and recall depends on the specific goals/context of the application.
 - ⇒ In some cases, such as detecting fraudulent transactions, precision may be more important than recall, as false positives can have significant consequences.
 - ⇒ In other cases, such as detecting rare medical conditions, recall may be more important, as missing a true positive can have serious consequences.

Compromises: Conclusions

1. There is always the need to compromise, or **tradeoff**:
 - (a) Bias-Variance
 - (b) Precision-Recall
2. The compromise we choose will always be subjective, and require **judgement in context**. There is no single, well-defined, “right” answer.
3. This will necessarily introduce an **ethical/cognitive bias**. See Ethics lectures.

MSE in Practice

See [Overview - general](#)

- This is the **oldest** (Gauss, Legendre, in the 1800's...) error (or loss) metric.
 - ⇒ MSE, or RMSE, is adapted to **optimization** errors (is this the best possible fit?), and should be used exclusively for this.
 - ⇒ MSE, or RMSE, is not adapted to **evaluation** of a ML model (how well is it doing?), since it is highly sensitive to outliers.
- For model performance evaluation, one should prefer the **Median Absolute Deviation** (MAD), or the Mean Absolute Error (MAE), though the former is more robust to outliers. In other words, the choice of metric depends on the characteristics of your data.
- To avoid orders of magnitude in errors, due to measurement units, and to obtain relative measures, the MAD and MAE can be recalculated as percentages.

R^2 in Practice

- This very widely-used regression metric has a number of shortcomings, and some people say that it should not be used at all, because:
 - ⇒ R^2 relies on data leakage
 - ⇒ R^2 decreases with increasing noise level
 - ⇒ R^2 increases with increasing model complexity
- One can use R^2 for the following instances:
 - ⇒ determine whether a change in explanatory variables improves the model fit
 - ⇒ compare models with different subsets of explanatory variables
 - ⇒ detecting co-linearity⁴
- Possible solutions

⁴High R-squared alongside insignificant coefficients can indicate multicollinearity, where independent variables are highly correlated and inflate the standard errors of the regression coefficients.

- ⇒ do not use it,
- ⇒ use an out-of-sample version where, for test performance evaluation, the TSS is calculated over the training samples only,

$$\text{TSS}_{\text{oos}} = \sum_{i=1}^n (y_i - \bar{y}_{\text{train}})^2$$

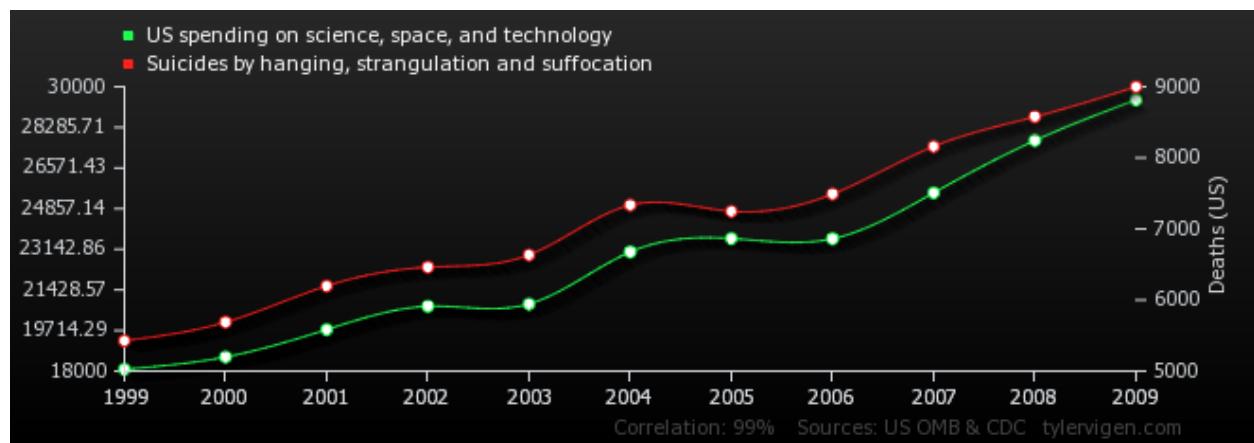
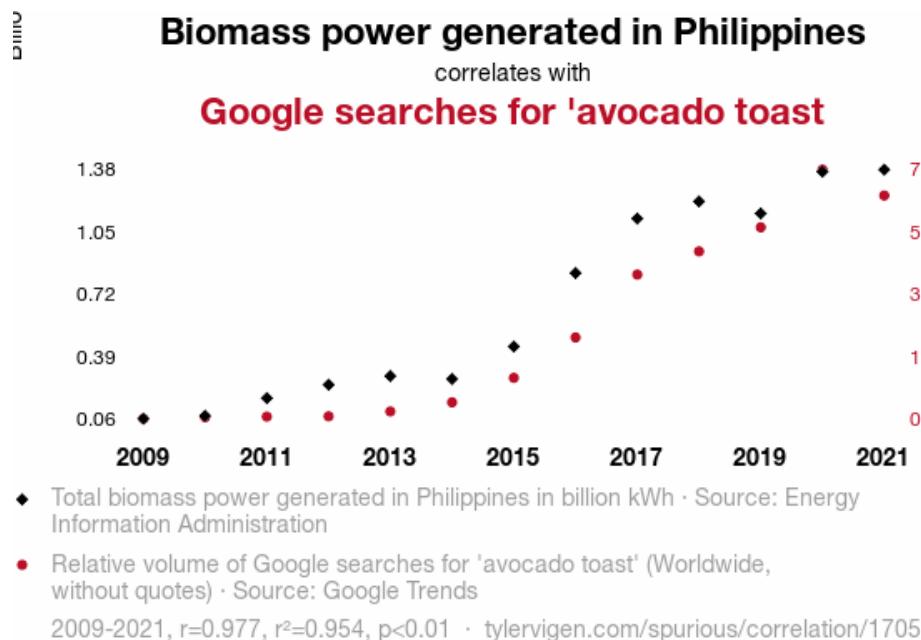
Classification Metrics

- There are a number of potential **pitfalls** in measuring the performance of classification models.
- The scikit [tutorial](#) is an excellent basis for illustrating several of these problems.
- See also: [Overview - general](#).

CAUSALITY and CORRELATION

Principles

- Correlation can be misleading.



Warning

Correlation does not always imply causation.

- Causality can be difficult to evaluate.

⇒ Confounder analysis.

⇒ Partial correlation.

⇒ Shapley values.

FEATURES

Features

- Feature **selection**: see [model selection](#).
- Feature **engineering**: transformed values can be more meaningful.
- Feature **importance**: use ranking provided by ensemble-based methods.
 - ⇒ The scikit [tutorial](#) is an excellent basis for illustrating these aspects.
 - ⇒ **Shapley** values are powerful indicators of feature importance.

Shapley Values

Definition

Shapley values quantify the significance of individual input features' contribution to the model's predicted values.

- Analysis of **explainability** through SHAP regression values aims to evaluate the contribution of input variables (often called “input features”) to the predictions made by a machine learning predictor model.
- The **contribution** of that input feature to a prediction is calculated mathematically through the construction of a so-called “explanation model”.
- The **explanation** model evaluates the predictions of a machine learning system as the sum of the contributions of each input feature and the mean predicted value.

- Mathematically, the explanation model can be stated as

$$y = \bar{y} + \sum_i \phi_i,$$

where

- ⇒ y is an individual prediction,
- ⇒ \bar{y} is the average predicted value across all predictions
- ⇒ ϕ_i is the contribution of input feature i to the prediction (also known as the “SHAP regression value” or “SHAP value”).
- Input variables with **larger** magnitude SHAP values are interpreted as contributing **more** to a specific prediction than those with a smaller magnitude SHAP values.
- For a given case,
 - ⇒ **positive** SHAP values indicate a specific feature contributes toward **increasing** the final predicted value y , and
 - ⇒ **negative** SHAP values indicate a contribution toward **decreasing** the prediction.

- These SHAP values, ϕ_i , are calculated following a **game theoretic** approach to assess prediction contributions.
- The SHAP framework has several key desirable properties:
 - ⇒ The **sum** of the contributions accurately reproduces the predicted value.
 - ⇒ The **contributions** of input features that are not present in the machine learning model are assigned values of **0**.
 - ⇒ SHAP is a machine learning **model-agnostic** technique and can be applied to a wide class of prediction methods.
- Despite these key advantages, there are several potential deficiencies to the application of SHAP analysis to error characterization in the Earth Sciences.
 - ⇒ SHAP analysis cannot directly yield **causal** insights.
 - ⇒ Direct calculation of SHAP values is very computationally **expensive**.
 - As such, care must be taken to properly interpret the SHAP values resulting from any particular analysis.

- It is important to note that to improve computational performance, common implementations of SHAP analysis are available in existing machine learning libraries.

Warnings on SHAP for Causality

- Flexible predictive models like XGBoost or LightGBM are powerful tools for solving **prediction** problems. However, they are **not inherently causal** models, so interpreting them with SHAP will fail to accurately answer causal questions in many common situations.
- Unless features in a model are the result of experimental variation, applying SHAP to predictive models without considering **confounding** is generally not an appropriate tool to measure causal impacts used to inform policy.
- SHAP and other interpretability tools can be useful for causal inference, and SHAP is integrated into many causal inference packages, but those use cases are **explicitly causal** in nature.
- To that end, using the same data we would collect for prediction problems and using **causal inference** methods like **double ML** that are particularly designed to return

causal effects is often a good approach for informing policy.

- In other situations, only an experiment or other source of **randomization** can really answer **what if** questions. Causal inference always requires us to make important assumptions.
- The main point is that the assumptions we make by interpreting a normal predictive model as causal, are often unrealistic.

PINN

Physics Informed Neural Networks

See: [Basics of NNs, PINNs](#)

Remark 4. This approach, though seemingly very elegant and powerful, has serious **problems** of cost, robustness and precision. It often requires a lot of tuning to work. PINNs should be employed with extreme attention and great caution.

ETHICS

Moral, Ethics and Bias

See: [Ethics](#)

References

1. Please consult the list provided on the CMU-IMU-2023 website:
[CODE REFERENCES](#)
2. Scikit-learn MOOC [website](#)
3. VSCODIUM for code editing [Site](#), [Downloads](#)
4. G. James, D. Witten, T. Hastie, R. Tibshirani, J. Taylor.
An Introduction to Statistical Learning. Springer. 2023.
[ISL Site](#)