

COMP90014

Algorithms for Bioinformatics

Week 11B: Model Selection | Tuning | Validation

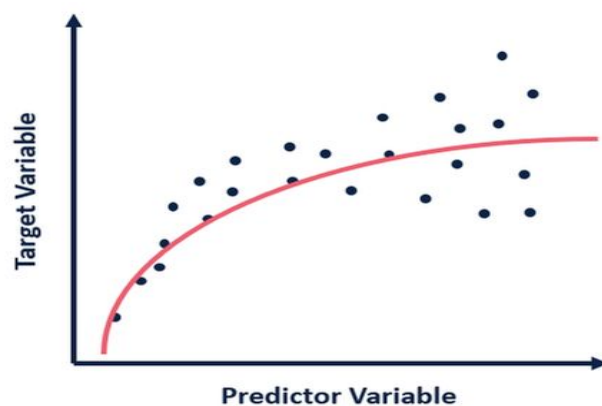
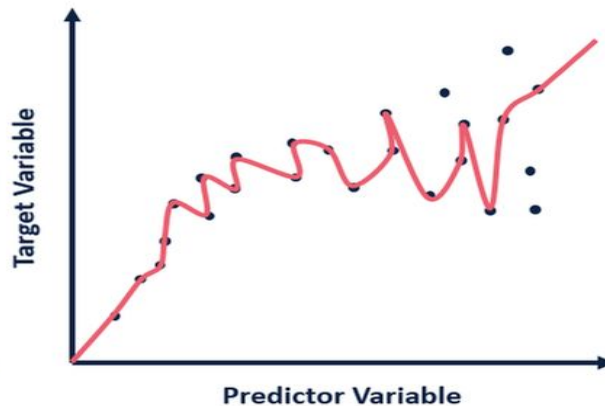
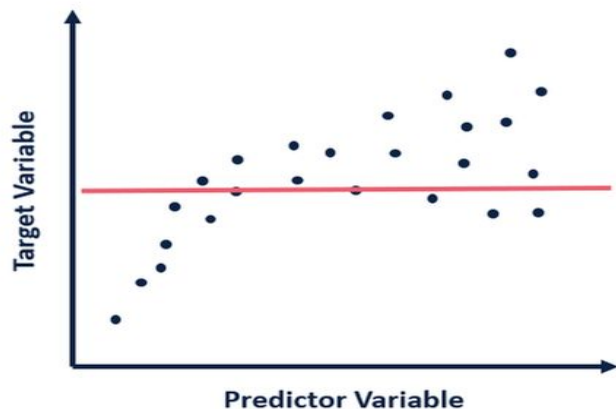


Supervised Learning

- Our goal: Generalization
- Underfitting and Overfitting
- **Learning Algorithms**
 - K-nearest neighbours (KNN)
 - Naïve Bayes
 - Decision Trees
 - Support Vector Machines (SVMs)
 - Ensemble methods
- **Model validation**
 - Hold-out & Cross-validation
- **Evaluation Metrics**
 - **Classification**
 - Confusion matrix
 - Type I/II errors
 - ROC curves
 - Imbalanced data
 - **Regression**
 - Correlation coefficient
 - Mean Squared Error
- **Tools and Packages**

Performance Estimation

- **Overfitting and underfitting**



- **Bias**

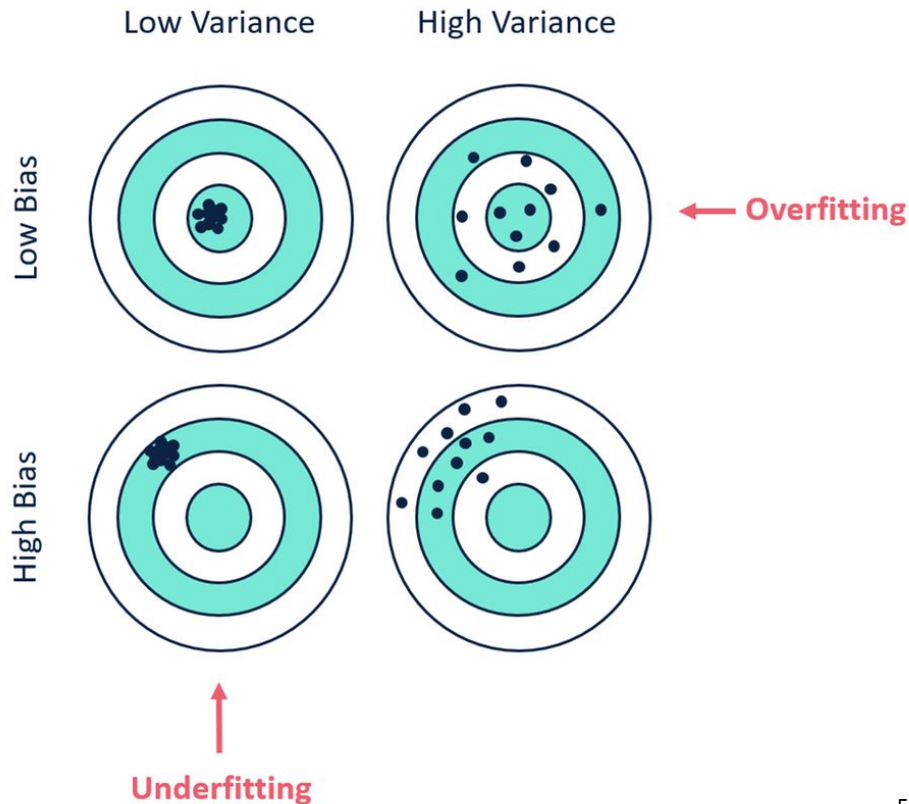
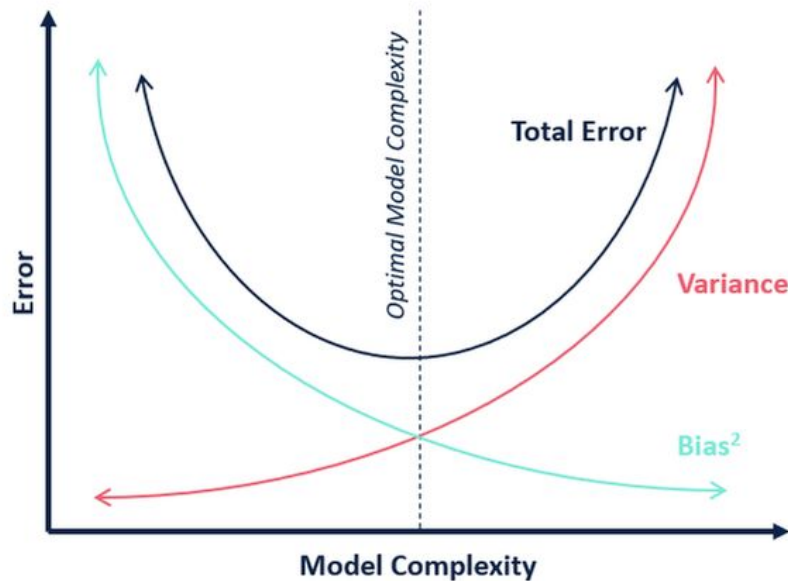
- Difference between prediction and real outcome

- **Variance**

- Variability of predictions

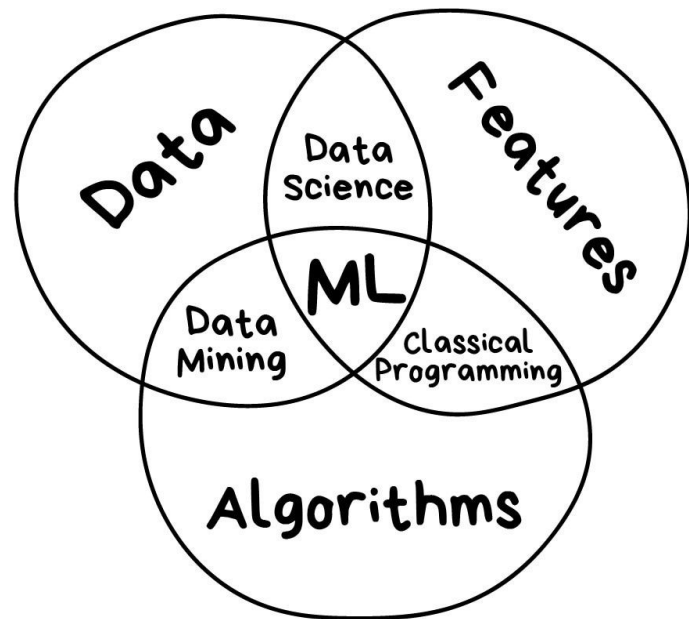
Performance Estimation

- **Overfitting**
 - Low bias, high variance
- **Underfitting**
 - High variance, low bias



Performance estimation

- How do we know that it **generalizes well to unseen data** rather than **simply memorize the training data**?
 - And how do we select a **good predictive model**?
 - Perhaps a different **algorithm** would be more appropriate?
 - Do we need more (or higher quality) **data**?
 - Do we need to investigate different **features**?

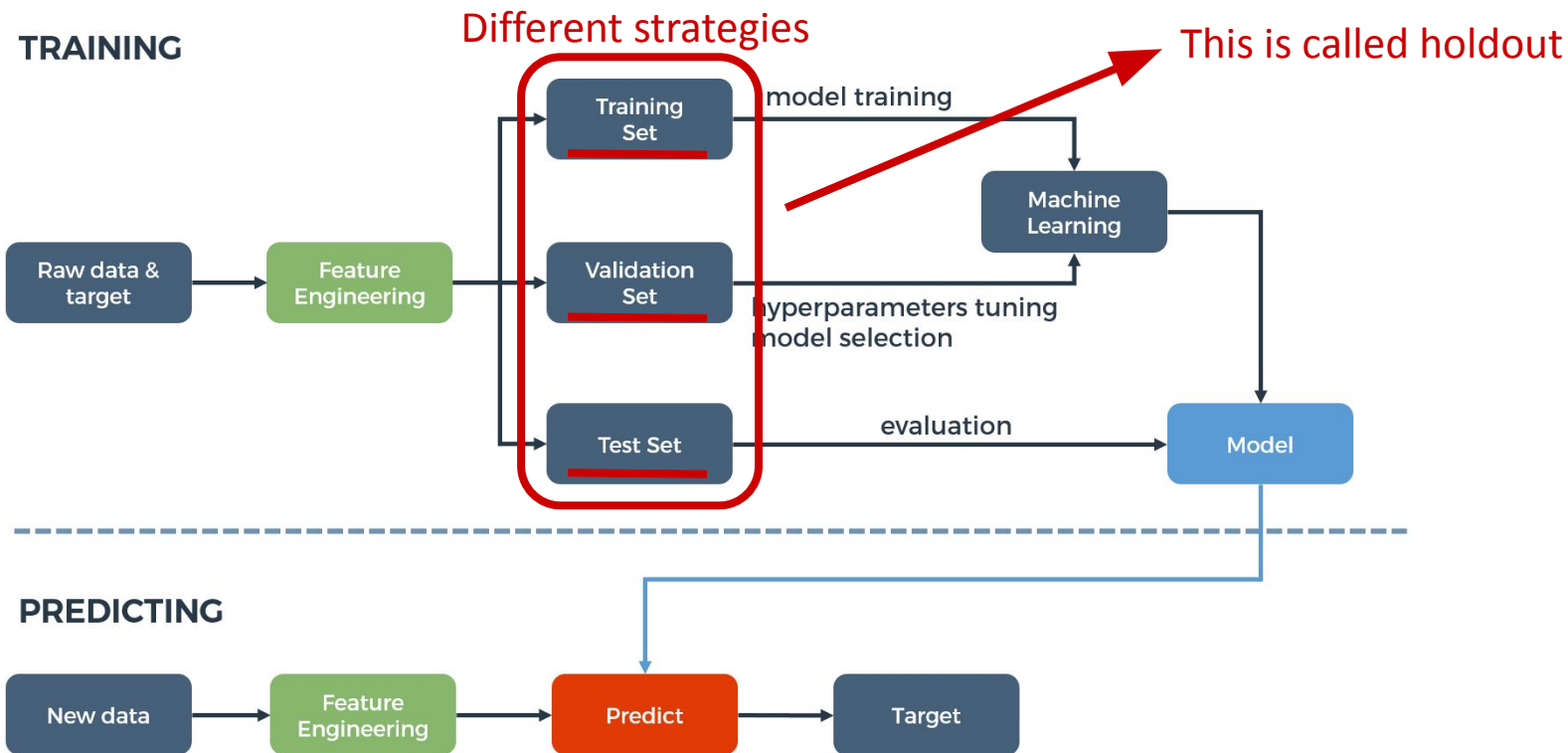


Performance Estimation

- Avoid **overfitting** and **underfitting**
 - A robust **validation strategy**
 - Choosing the right **performance metrics**
- We want to estimate the **generalization performance**
 - Predictive performance on **unseen data**
- We need to be able to **compare predictive models** to choose the **best/most appropriate one**
 - Assessing different **algorithms**,
 - **Parameters** and
 - **Feature** combinations



Model Selection and Validation

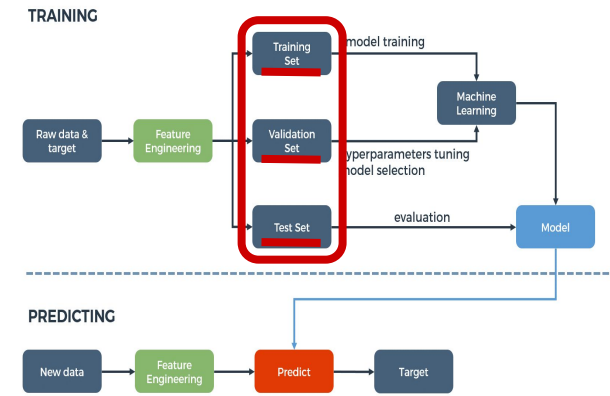




Model Validation: Holdout

- **Holdout Validation**

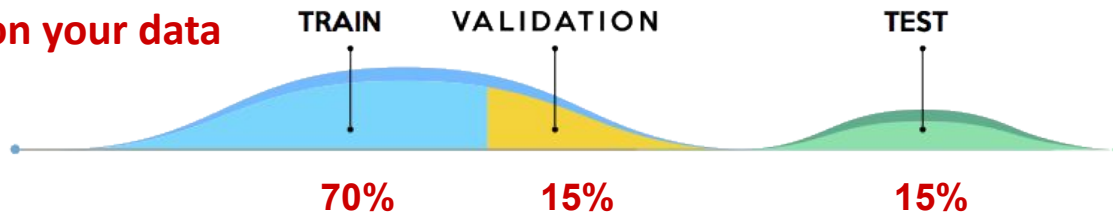
- You can't use the same data to train and test the model
 - That would be cheating!
- Test a model on different data than it was trained on
 - Provides an unbiased estimate of learning performance
- **Holdout:** dataset is randomly divided into three subsets
 - Training set to build predictive models
 - Validation set to assess performance in training and to fine-tune parameters
 - Independent test set (blind test) to assess the likely performance on unseen data
- What if performance on training is much better than on the test set? **Overfitting**



HOLDOUT STRATEGY

Highly depends on your data

1 Split your data into train / validation / test



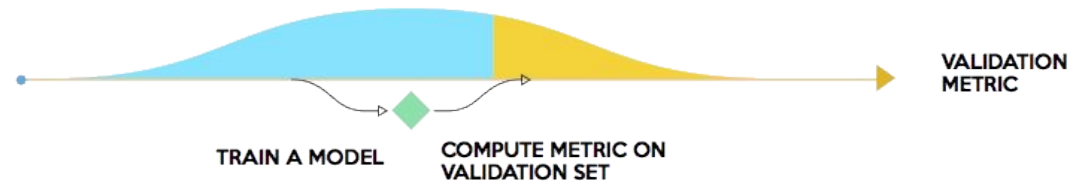
2 For each parameter combination

Parameter A (e.g., depth)

4	14
5	15
6	16
7	17

 Parameter B (e.g., n trees)

4	14
5	15
6	16
7	17



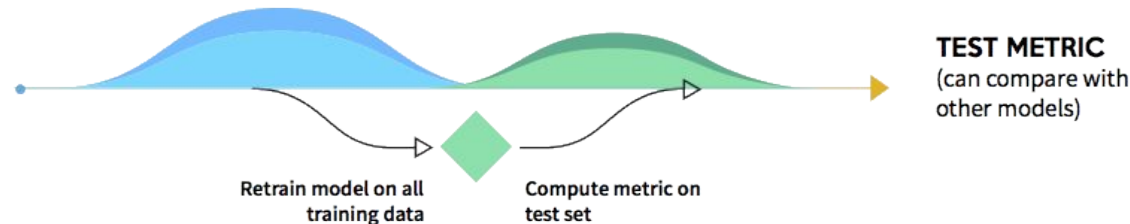
3 Choose the parameter combination with the best metric

Parameter A (e.g., depth)

6	14
---	----

 Parameter B (e.g., n trees)

4	14
---	----





Model Evaluation

- **Holdout Strategy Limitation**
 - Prone to **selection bias**
- **Solution**
 - **Repeated** Holdout Validation
 - 100x, 1000x
 - **Bootstrapping**
 - Resampling with replacement

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------

x_8	x_6	x_2	x_9	x_5	x_8	x_1	x_4	x_8	x_2
-------	-------	-------	-------	-------	-------	-------	-------	-------	-------

x_3	x_7	x_{10}
-------	-------	----------

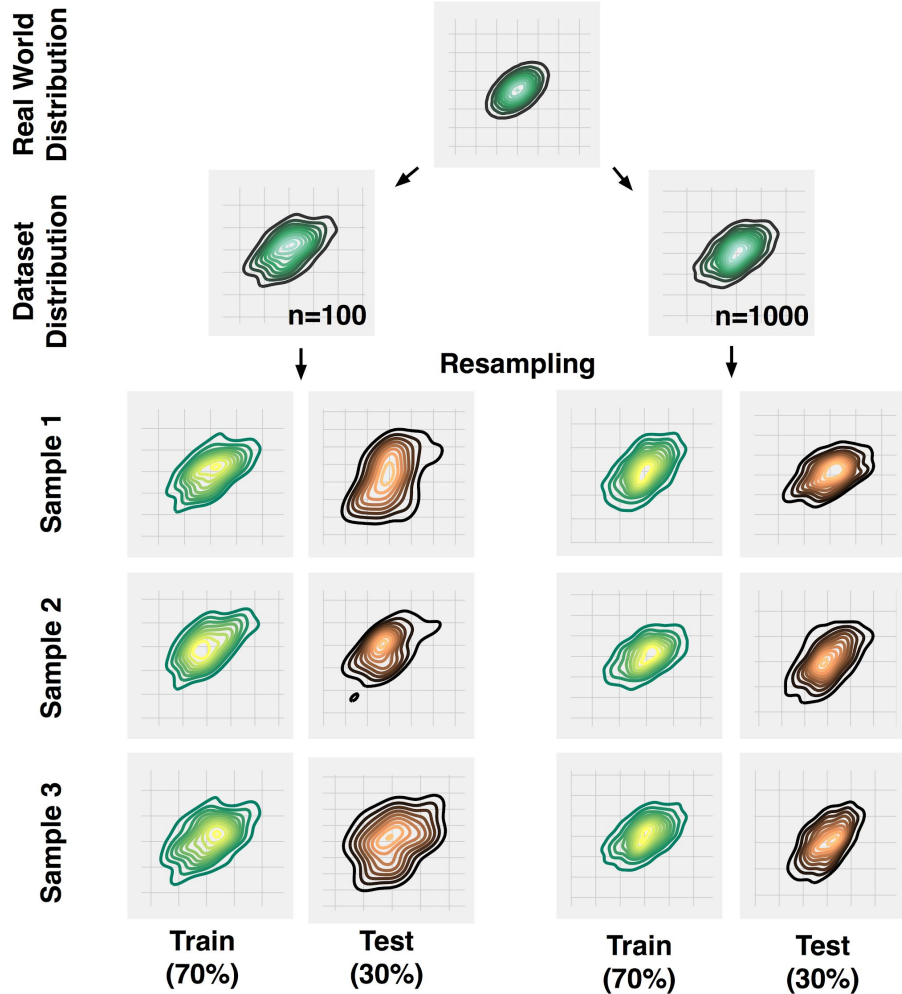
x_{10}	x_1	x_3	x_5	x_1	x_7	x_4	x_2	x_1	x_8
----------	-------	-------	-------	-------	-------	-------	-------	-------	-------

x_6	x_9
-------	-------

x_6	x_5	x_4	x_1	x_2	x_4	x_2	x_6	x_9	x_2
-------	-------	-------	-------	-------	-------	-------	-------	-------	-------

x_3	x_7	x_8	x_{10}
-------	-------	-------	----------

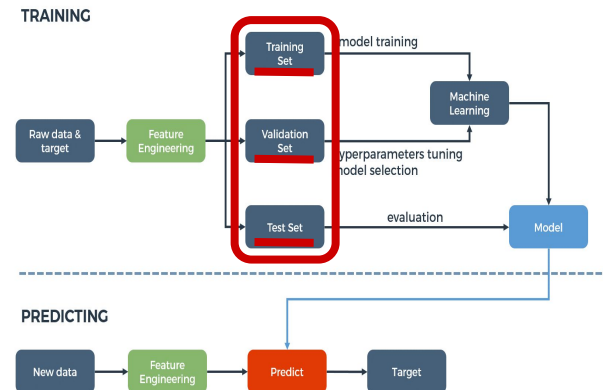
Test Sets



Model Validation: K-fold Cross Validation

- **K-fold Cross Validation**

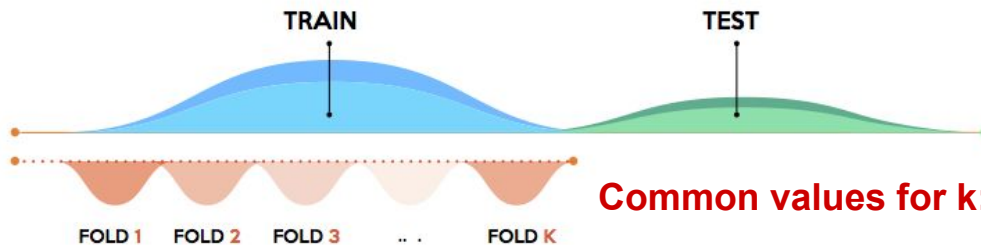
- Partition the dataset into k equal size subsamples (folds)
- We will iteratively use one of the subsets as validation and the other k-1 subsets are put together to form a training set
- Performance estimation is averaged over all k iterations
- Every data point gets to be in a validation set exactly once and gets to be in a training set k-1 times
- Reduces bias
 - Most of the data for training
- Reduces variance
 - All data is eventually used for validation



- We still need a independent test set

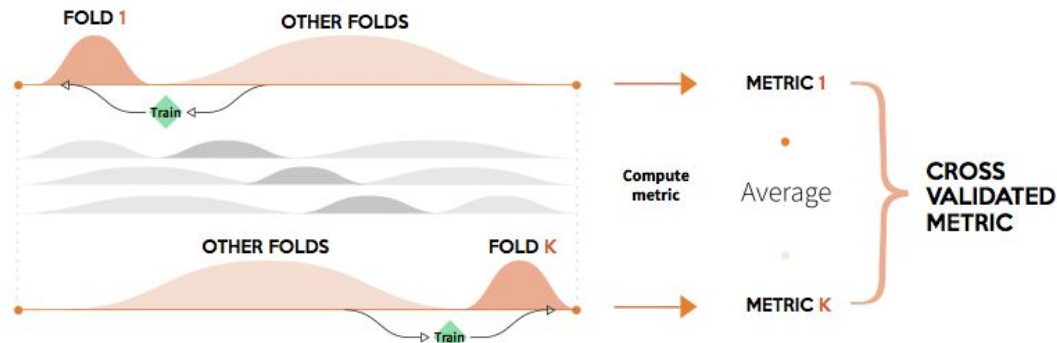
K-FOLD STRATEGY

1 Set aside the test set and split the train set into k folds

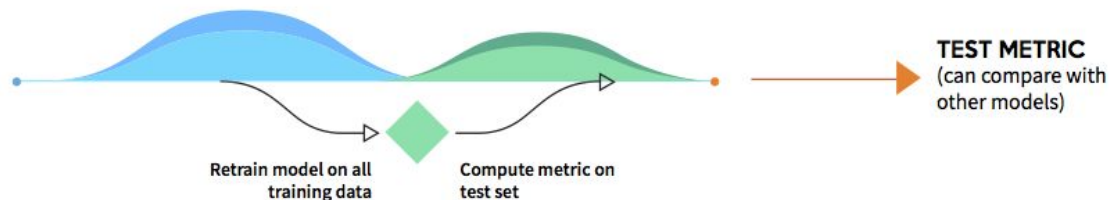


Common values for k: 5, 10, 20

2 For each parameter combination

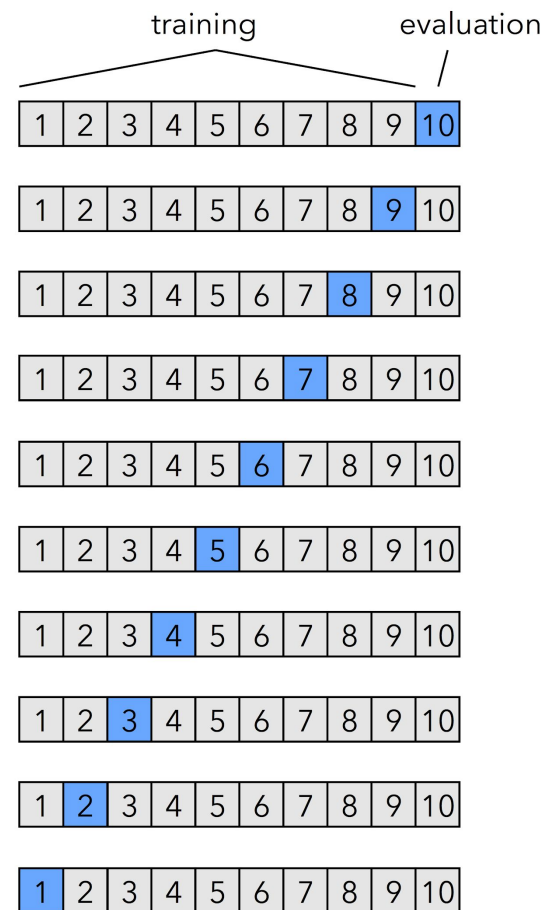


3 Choose the parameter combination with the best metrics



Cross Validation

- **Leave-one-out Cross Validation (LOOCV)**
 - K-fold Cross Validation when **K equals the number of data points we have for training/validation**
 - Per round, **each data point is considered individually in the validation/evaluation set**
 - **Maximizes** the amount of information available for **training**
 - Good for **small data sets**
 - **Computationally intensive**



Predictive Performance Metrics

- To select the **best performing model** we need to be able to **compare them**
- **Predictive performance metrics**
 - On the **validation** set
 - On an **independent test set**
 - Make sure they are **consistent**
- There are multiple metrics for **classification** and **regression**
- **Classification**
 - We can derive several metrics from a **confusion matrix**

Statistical Classification Metrics

<div><div>Sensitivity Recall Power</div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div>True Positive Rate</div></div>	TP	FP	FN	TN	TP	FP	FN	TN	<div><div>Precision</div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div>Positive Predictive Value</div></div>	TP	FP	FN	TN	TP	FP	FN	TN	<div><div></div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div>False Discovery Rate</div></div>	TP	FP	FN	TN	TP	FP	FN	TN	<div><div>Type I Error α Fall Out</div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div>False Positive Rate</div></div>	TP	FP	FN	TN	TP	FP	FN	TN	<div><div>Accuracy</div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div></div>	TP	FP	FN	TN	TP	FP	FN	TN	<div><div>F1 Score F Measure</div><div><table><tr><td>2x TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div><table><tr><td>2x TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div></div>	2x TP	FP	FN	TN	2x TP	FP	FN	TN							
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
2x TP	FP																																																											
FN	TN																																																											
2x TP	FP																																																											
FN	TN																																																											
<div><div>Type II Error β</div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div>False Negative Rate</div></div>	TP	FP	FN	TN	TP	FP	FN	TN	<div><div></div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div>True Discovery Rate</div></div>	TP	FP	FN	TN	TP	FP	FN	TN	<div><div></div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div>Negative Predictive Value</div></div>	TP	FP	FN	TN	TP	FP	FN	TN	<div><div>Specificity</div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div>True Negative Rate</div></div>	TP	FP	FN	TN	TP	FP	FN	TN	<div><div>Confusion Matrix</div><div><table><tr><td></td><td></td><td colspan="2">actual</td></tr><tr><td></td><td></td><td>T</td><td>F</td></tr><tr><td rowspan="2">predicted</td><td>P</td><td>TP</td><td>FP</td></tr><tr><td>N</td><td>FN</td><td>TN</td></tr></table></div><div>TP: True Positive FP: False Positive FN: False Negative TN: True Negative</div><div>actual = observed predicted = expected</div></div>			actual				T	F	predicted	P	TP	FP	N	FN	TN	<div><div>Matthews Correlation Coefficient</div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div><table><tr><td>TP</td><td>FP</td></tr><tr><td>FN</td><td>TN</td></tr></table></div><div>difference of products</div><div>square root of product of sums</div></div>	TP	FP	FN	TN	TP	FP	FN	TN
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											
		actual																																																										
		T	F																																																									
predicted	P	TP	FP																																																									
	N	FN	TN																																																									
TP	FP																																																											
FN	TN																																																											
TP	FP																																																											
FN	TN																																																											

Confusion Matrix

- **Table layout** that allows **visualisation** of predictive performance
 - **Rows** represent the instances in a **predicted class**
 - **Columns** represent the instances in an **actual class**
- For binary classification (two classes, **positive** and **negative**)
- **True Positives (TP)**
 - **Correctly predicted** as belonging to the **positive class**
 - A cancer test correctly identifying a patient who has cancer
- **True Negatives (TN)**
 - **Correctly predicted** as belonging to the **negative class**
 - A cancer test correctly identifying a patient who doesn't have cancer

		actual	
		T	F
predicted	P	TP	FP
	N	FN	TN

TP: True Positive
FP: False Positive
FN: False Negative
TN: True Negative

Confusion Matrix

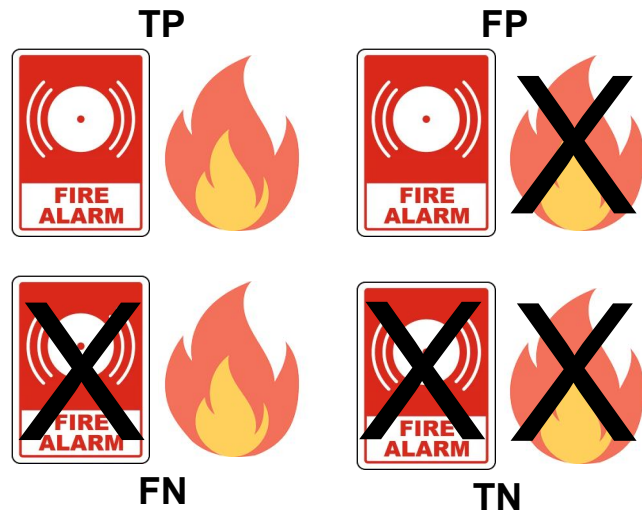
- **False Positives (FP) - Type I Error**
 - **Incorrectly predicted** as belonging to the **positive class**
 - A cancer test saying a patient has cancer, while they actually don't
- **False Negatives (FN) - Type II Error**
 - **Incorrectly predicted** as belonging to the **negative class**
 - A cancer test saying a patient doesn't have cancer, while they actually do
- **We want to:**
 - **Minimize False*** & **Maximize True*** cases
 - **Which Type Error is worse?**
 - Depends on the problem

		actual	
		T	F
predicted	P	TP	FP
	N	FN	TN

TP: True Positive
FP: False Positive
FN: False Negative
TN: True Negative

Type I and II errors

- Which Type Error is worse?
 - Fire alarm
- Type I:
 - Fire alarm **rings** when there is **no fire**
- Type II:
 - Fire alarm **fails to ring** when there **is fire**
- Which Type Error is worse in this case?



Accuracy

- **Accuracy** is the number of **correct predictions** made by the model **over all predictions** made
 - $Accuracy = (TP+TN) / (TP+TN+FP+FN)$
- **When can I use Accuracy?**
 - Accuracy is a good measure when the **target classes in the data are nearly balanced**
 - Similar number of data points belonging to each class
- **When NOT to use Accuracy?**
 - **Imbalanced data sets**
 - *e.g.*, 95% of the data belong to class A, 5% to class B
 - A predictor that **only guesses class A** has 95% accuracy

Accuracy

TP	FP
FN	TN

TP	FP
FN	TN

Accuracy

Spam filter (25 spam messages, 125 not spam)

- 73.3%

	Spam	Not spam
Pred. spam	10 (TP)	25 (FP)
Pred. not-spam	15 (FN)	100 (TN)

- 83.3%

	Spam	Not spam
Pred. spam	0 (TP)	25 (FP)
Pred. not-spam	0 (FN)	125 (TN)

Accuracy

TP	FP
FN	TN

TP	FP
FN	TN

Precision

- **Precision** is the proportion of **predicted positives** that truly are **positives**
 - $\text{Precision} = (TP) / (TP + FP)$
- Takes into account **Type I Error**
- **Precision** is a valid choice of evaluation metric when we **want to be very sure** of our **positive class prediction**.
 - *e.g.*, if we want to **predict if we should decrease the credit limit on a particular account**
 - We want **minimum FP** otherwise it may result in customer dissatisfaction
 - Maximise **precision**

Precision

TP	FP
FN	TN

TP	FP
FN	TN

Precision

Spam filter (10 spam messages, 90 not spam)

	Spam	Not spam
Pred. spam	1 (TP)	0 (FP)
Pred. not-spam	9 (FN)	90 (TN)

- What's the **precision** for the spam class?
 - **100%**

Precision

TP	FP
FN	TN

TP	FP
FN	TN

Recall or Sensitivity

- **Recall** is the proportion of **actual positives** that are **correctly classified**
 - $Recall = (TP) / (TP + FN)$
- Takes into account **Type II Error**
- **Recall** is used when we **want to recover as many positives as we can**
 - e.g., If we want to **predict if a patient has a disease or not**, we want to capture the disease even if we are not very sure
 - We want **minimum FN** otherwise we might discharge a patient that needs treatment - maximise **recall**
- **Caveat**
 - If we predict everything as positive, **recall will be 100%**

Sensitivity Recall Power

TP	FP
FN	TN

TP	FP
FN	TN

Recall or Sensitivity

Spam filter (10 spam messages, 90 not spam)

	Spam	Not spam
Pred. spam	10 (TP)	90 (FP)
Pred. not-spam	0 (FN)	0 (TN)

- What's the **recall** for the spam class?
 - **100%**

**Sensitivity
Recall
Power**

TP	FP
FN	TN

TP	FP
FN	TN

F1-score or F1-measure

- How to find a compromise between precision and recall?
 - Does simply taking their **arithmetic mean** work?
 - *e.g.*, a predictive model with **20% recall** and **100% precision**
 - `mean(Precision+Recall) = 60%`
- **F1-score** is the harmonic mean of precision and recall
 - `F1 = 2*(Precision*Recall)/(Precision+Recall)`
 - For the example: `F1 = 33%`
- F1-score will **penalise large discrepancies between precision and recall**

F1 Score F Measure

2x TP	FP
FN	TN

2x TP	FP
FN	TN

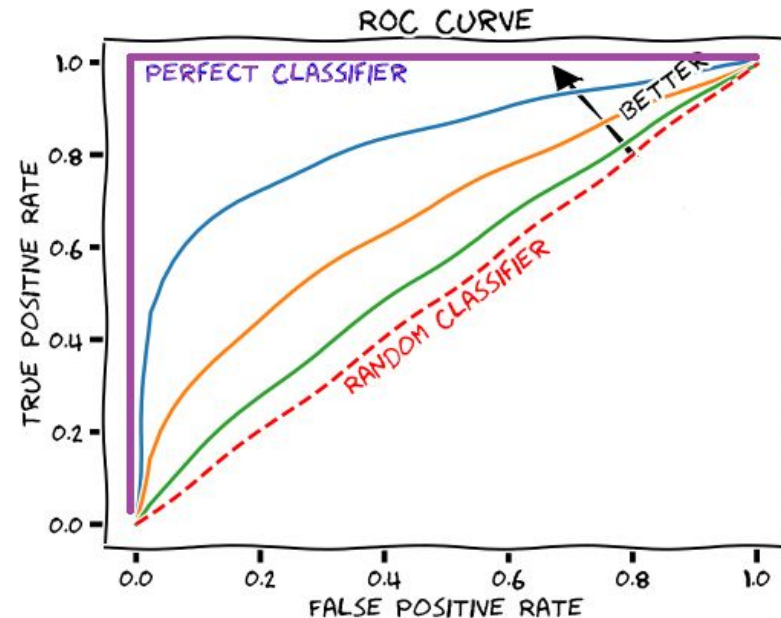
Receiver Operating Characteristic Curves

- **ROC curves**

- Is a graph showing the performance of a classification model at **different thresholds**
 - *i.e.*, **class probabilities** from our classifier
- Y-axis: **true positive rate (recall)**
- X-axis: **false positive rate (1-specificity)**
 - $\text{Specificity} = \text{TN}/(\text{TN} + \text{FP})$

- **AUC (Area Under the ROC Curve)**

- Varies from 0 to 1
- A **random binary classifier**: AUC of 0.5



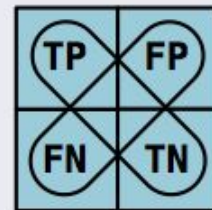
Check it out! (StatQuest Channel)

<https://www.youtube.com/watch?v=4jRBRDbJemM>

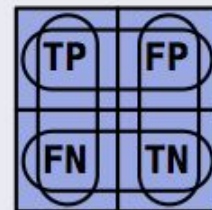
Matthews Correlation Coefficient (MCC)

- **Matthews Correlation Coefficient (MCC)** takes into account **true and false positives and negatives**
 - It is considered a **balanced metric**
- Very good metric for **imbalanced** data sets
 - Even classes are of **very different in sizes**
 - In contrast with **accuracy**
- **Ranges between -1 and 1**
 - 1 shows a **perfect prediction**
 - 0 equals to the **random prediction**
 - -1 indicates **total disagreement between predicted and actual labels**

Matthews Correlation Coefficient



difference of products



square root of product of sums

$$MCC = \frac{TP \times TN - FN \times FP}{\sqrt{(TP + FN)(TP + FP)(TN + FN)(TN + FP)}}$$

Evaluating regression models

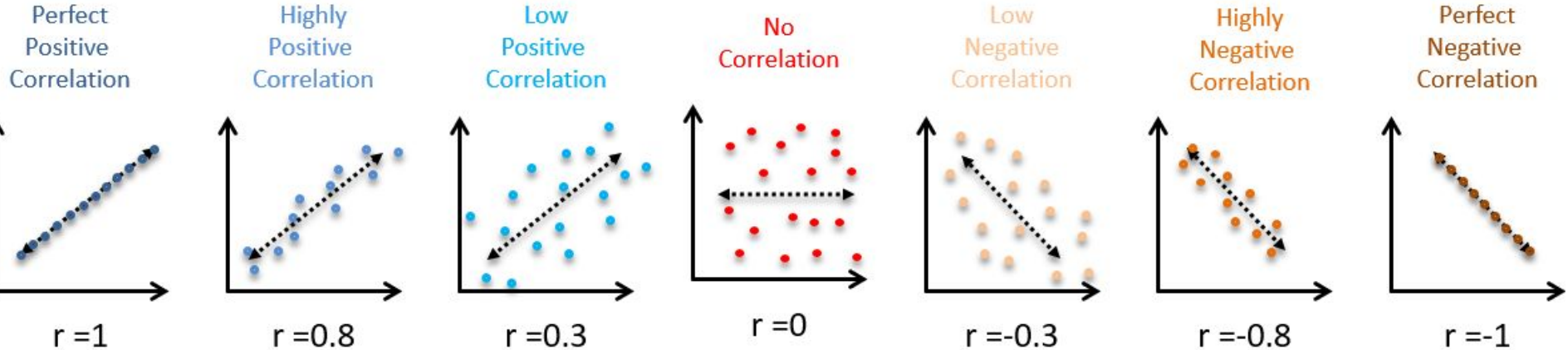
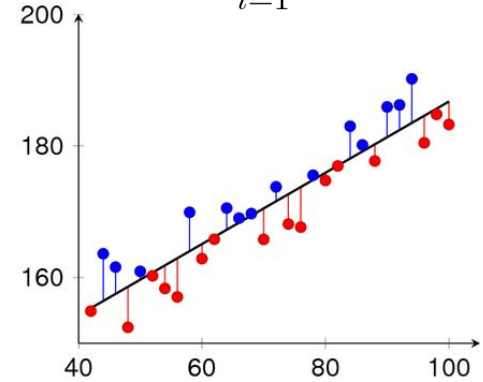
- **Mean Square Error (MSE)**

- The average of squared differences between the predicted and the actual values

- **Pearson Correlation Coefficient**

- A measure of the linear correlation between two variables (predicted vs. actual)

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



Which models is the best one?

Metric	MODEL1	MODEL2
Recall	0.6667	0.8333
Specificity	0.8333	0.6667
Precision	0.8000	0.7143
Accuracy	0.7500	0.7500
F1 Score	0.7273	0.7692
MCC	0.5071	0.5071

- **Model 1**
 - To minimise **Type I Error (better precision)**

- **Model 2**
 - To minimise **Type 2 Error (better recall)**

- **Benign vs. malignant cancer hypothetical cases**

MODEL 1	Actual disease	Actual healthy
Predicted disease	200	50
Predicted healthy	100	250

MODEL 2	Actual disease	Actual healthy
Predicted disease	250	100
Predicted healthy	50	200

Tools and Packages



- **Scikit Learn**

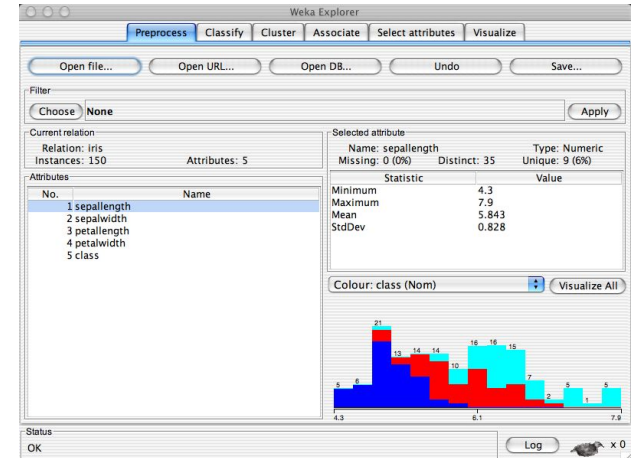
- Machine Learning in **Python**
- Data analysis
- Built on NumPy, SciPy, and matplotlib
- Open source

- **TensorFlow**

- **Python**
- Developed by Google
- **Deep Learning**

- **Weka**

- GUI
- Java



Thank you!

Today: Model selection, tuning, validation

Next time: Recap I