



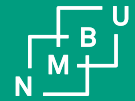
# Scikit-learn and Tour of Classifiers

Recap



## Classification pipeline

- **Goal:** Train a classifier that **generalizes** well on **unseen** data
- Train (for training) / Test (unseen) split to get a better estimate of the generalization error via the error on the test set
  - **Randomize and stratify** to get equal class distributions in the test and training set
- Feature Scaling/Transformations
  - Always apply **identical** transformations to the test and training set
  - **Avoid information leakage!** Standardization: Compute mean/std **on the training set**
- **Accuracy:**  $1.0 - \text{misclassified\_samples} / \text{all\_samples}$



# Overfitting / Underfitting



## Overfitting

- Overfitting is a common problem in machine learning
- We want the model to **generalize** well to **unseen** data
- A model that is overfitted performs well on the training data but does not generalise well to unseen data (test data)



## Overfitting – Underfitting

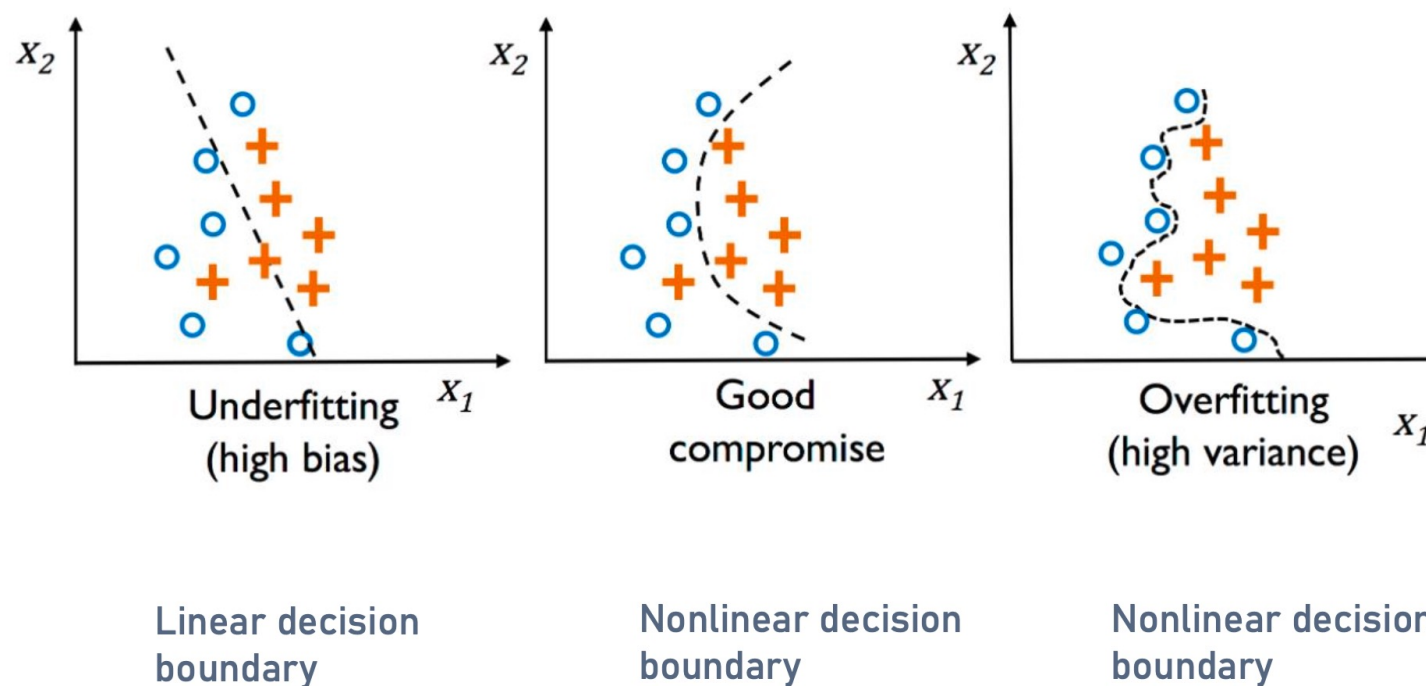
A model that suffers from **overfitting**

- Has high variance
- High variance may be caused by too many parameters that lead to a model that is too complex given the training data
- **poor performance** on unseen data because of high variance

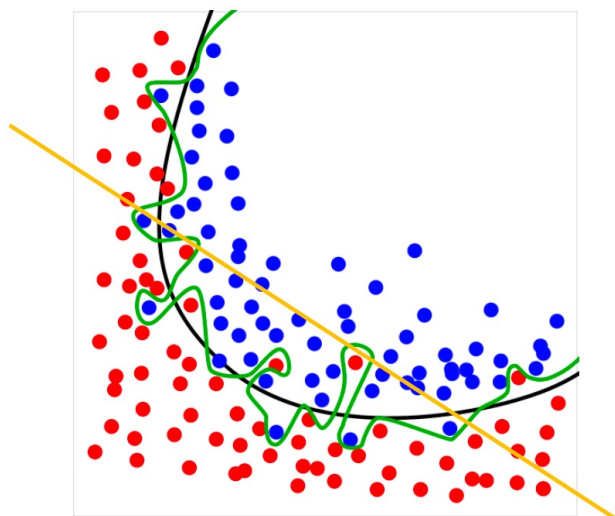
A model that suffers from **underfitting**

- Has high bias
- High bias means that the model is not complex enough to capture the patterns in the training data well
- **poor performance** on unseen data because of high bias

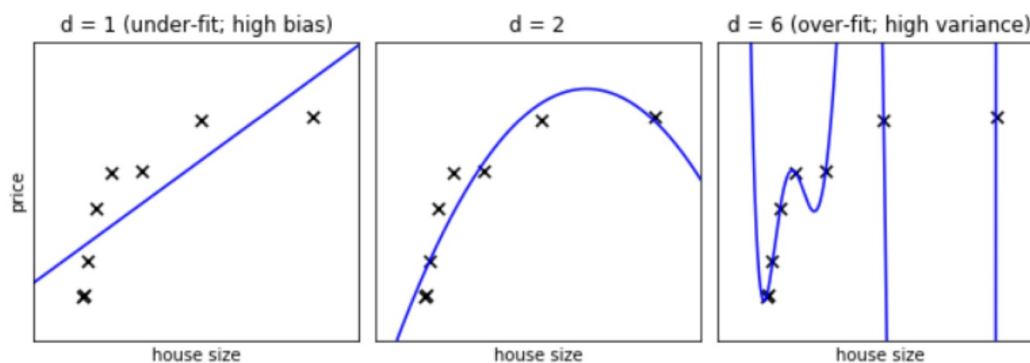
# Overfitting – Underfitting (Examples)



# Overfitting – Underfitting (Examples)

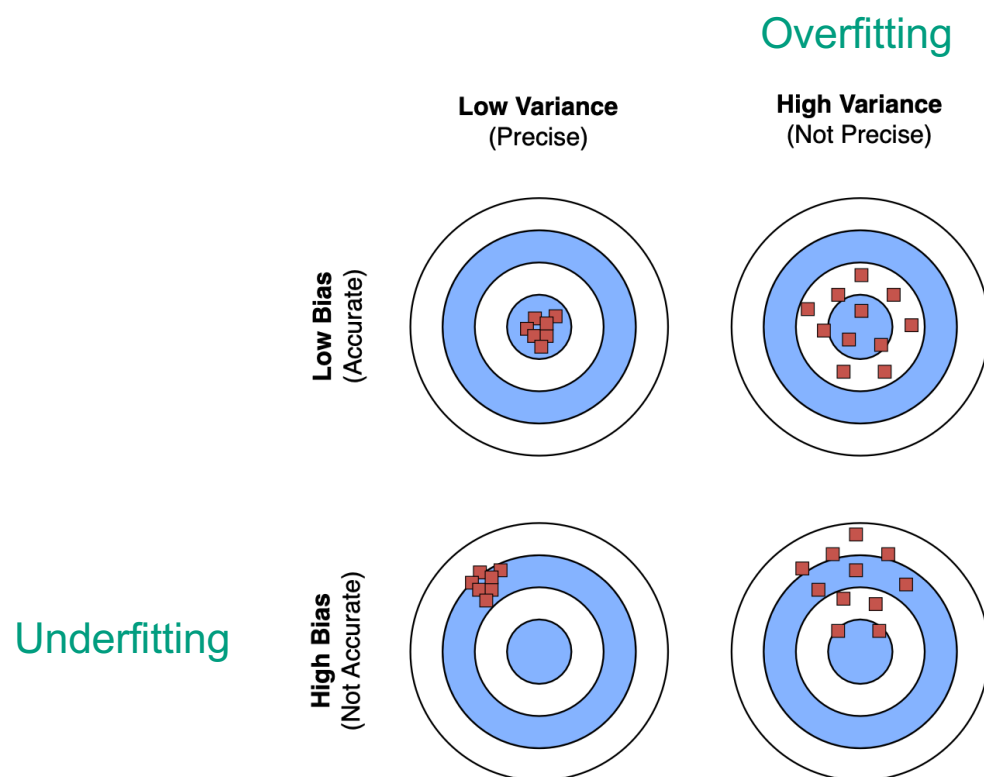


Classification



Regression

# Bias – Variance and Overfitting – Underfitting







## Bias – Variance

### Variance

- Measures the **consistency (variability)** of the model **prediction** for a particular sample
- If the model is retrained on different subsets of the training data and the prediction for a particular sample differs a lot each time → **high variance**
- The model is **sensitive** to the **randomness** in the **training data**

### Bias

- Measures **how far off** the predictions are from the **correct values** in general
- If the model is retrained multiple times on different **training** datasets, bias measures the **systematic error** that is **not due** to randomness



## Bias – Variance

Given a true value  $y$  and an estimator  $\hat{y}$

### Bias

$$Bias(\hat{y}) := E[\hat{y}] - y$$

### Variance

$$Var(\hat{y}) := E[(E[\hat{y}] - \hat{y})^2] = E[E[\hat{y}]^2 + \hat{y}^2 - 2E[\hat{y}]\hat{y}]$$

$$= E[\hat{y}]^2 + E[\hat{y}^2] - 2E[\hat{y}]E[\hat{y}] = E[\hat{y}]^2 + E[\hat{y}^2] - 2E[\hat{y}]^2 = E[\hat{y}^2] - E[\hat{y}]^2$$

$E[\hat{y}]$  : The expected value (expectation) of the estimator  $\hat{y}$  (here: expectation over the training sets)



## Bias – Variance decomposition of the squared loss

Given a true value  $y$  and an estimator  $\hat{y}$ ;  $Bias(\hat{y}) = E[\hat{y}] - y$ ;  $Var(\hat{y}) = E[(E[\hat{y}] - \hat{y})^2]$

For  $y$  (true outcome) and  $\hat{y}$  (estimated outcome):

Squared loss:  $L = (y - \hat{y})^2 = (y - E[\hat{y}] + E[\hat{y}] - \hat{y})^2$

Using:  $(a + b)^2 = a^2 + b^2 + 2ab$

$$= (y - E[\hat{y}])^2 + (E[\hat{y}] - \hat{y})^2 + 2(y - E[\hat{y}])(E[\hat{y}] - \hat{y})$$

Take the expected value of both sides:

$$E[L] = E[(y - E[\hat{y}])^2 + (E[\hat{y}] - \hat{y})^2] = Bias^2 + Variance$$

Using:  $E[2(y - E[\hat{y}])] = 2(y - E[\hat{y}])$  and  $E[(E[\hat{y}] - \hat{y})] = (E[\hat{y}] - E[\hat{y}]) = 0$



## Bias – Variance trade-off

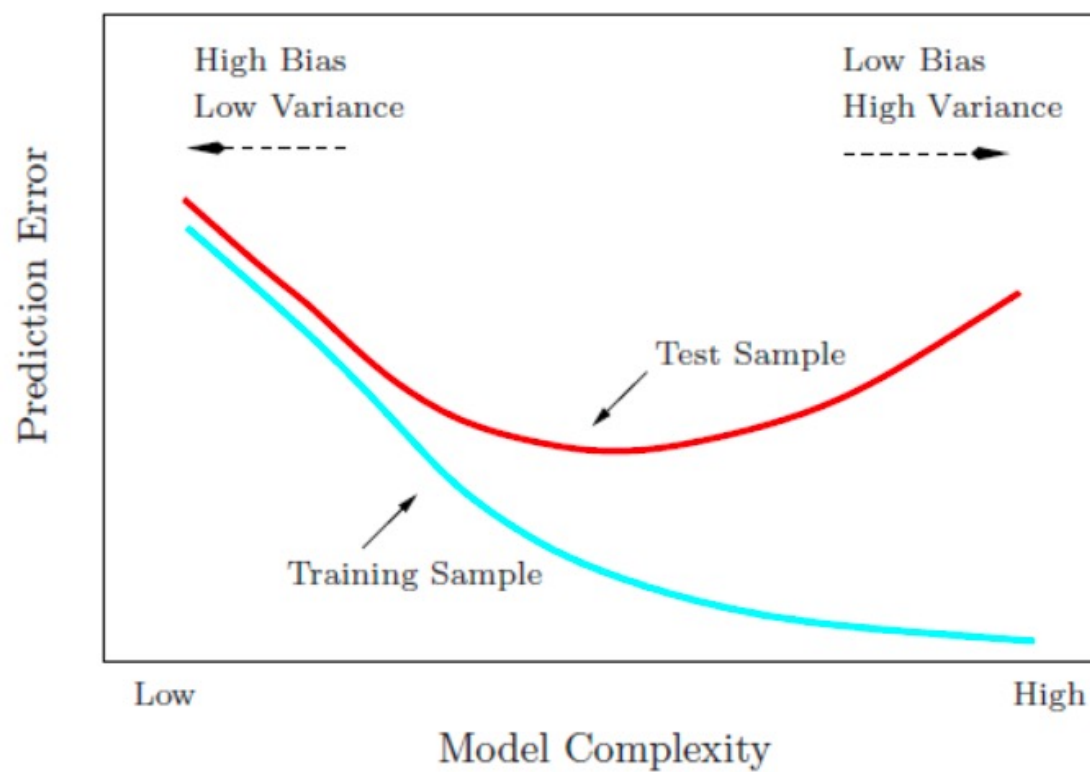


Figure from Elements of Statistical Learning (Hastie, Tibshirani, Friedman)



## Bias – Variance trade-off

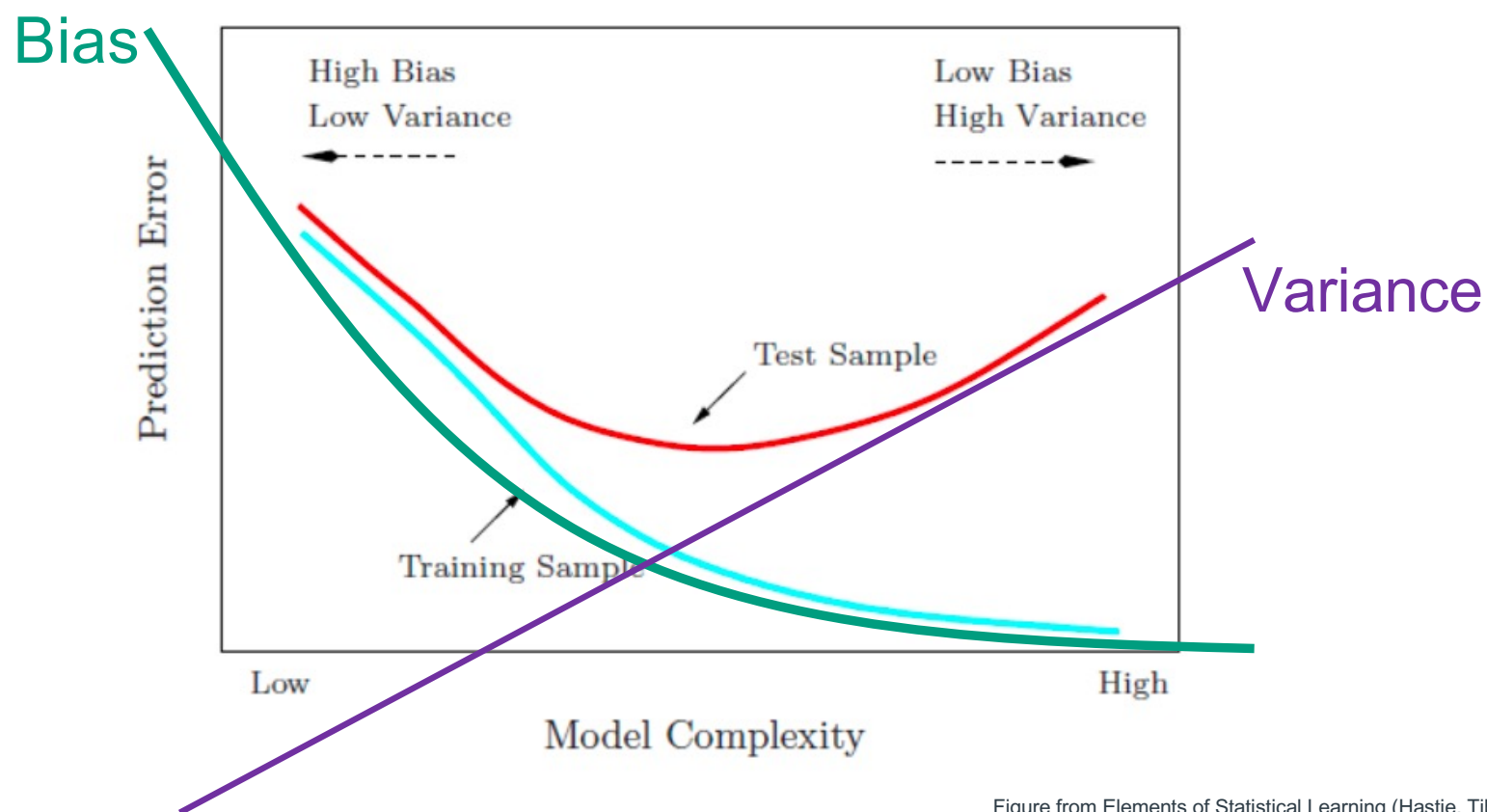


Figure from Elements of Statistical Learning (Hastie, Tibshirani, Friedman)



## Bias – Variance trade-off

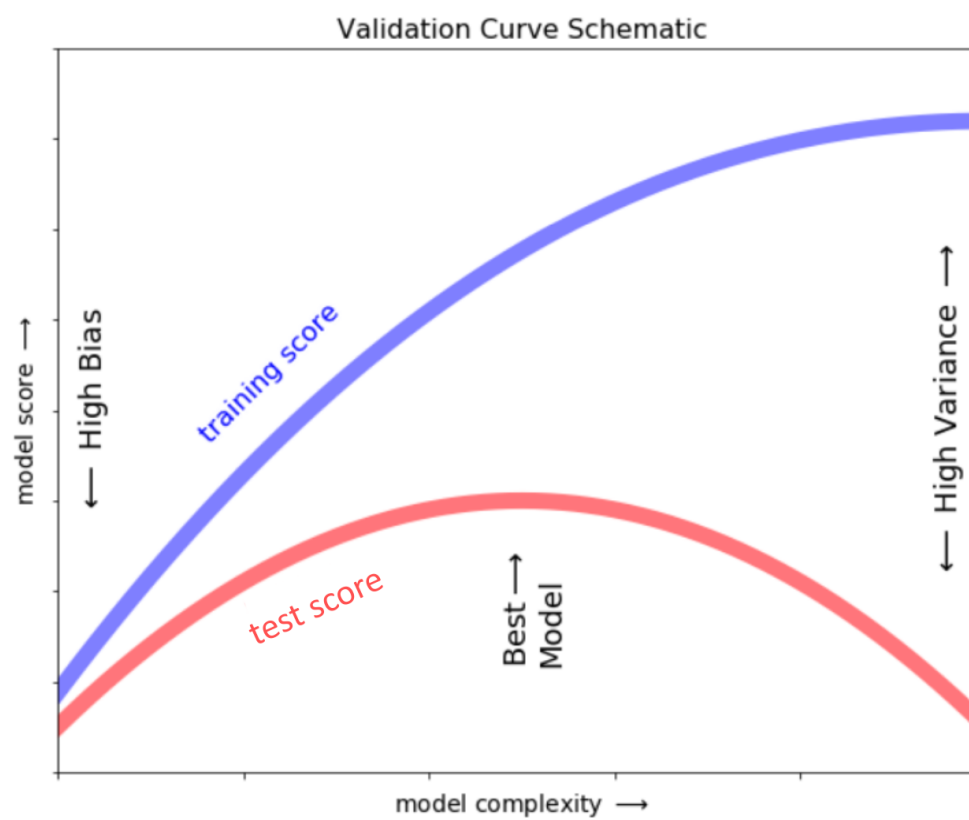


Figure from Python Machine Learning (Raschka & Mirjalili)



# Tackling overfitting via regularization



## Tackling overfitting

- Finding the appropriate **complexity – bias-variance trade-off** is important for good performance
  - Collect **more training data**
  - Introduce a **penalty for complexity** via **regularization**
  - Choose a **simpler model** with fewer parameters
  - **Reduce the dimensionality** of the data (feature selection, dimensionality reduction, Ch.06)





## Tackling overfitting via regularization

- Good option: **Tune** the complexity of the model using **regularisation**
- Regularisation is very useful for
  - Handling **collinearity** (high **correlation** among **features**)
  - **Filtering out noise** from the data
  - **Preventing overfitting**
- To make regularisation work properly, features need to be **scaled** / **standardised**
- **Concept** of regularisation:
  - Introduce **additional information (bias)** to penalise **extreme parameter** (weight) **values**
  - The **most common** form of regularisation is so-called **L2 regularisation** (sometimes also called **L2 shrinkage** or **weight decay**) (For details: See Ch.04 in the course book)



## $\ell_2$ - regularization (or L2 regularization / Ridge)

- $\ell_2$ -regularization: Add the Euclidean norm (two-norm) of the weights to the loss function
- $\ell_2$ -regularized logistic regression:

$$L(\boldsymbol{\theta}) = - \sum_{i=1}^n \left[ y^{(i)} \log(\sigma(z)) + (1 - y^{(i)}) \log(1 - \sigma(z)) \right] + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

Recall that  $\|\mathbf{w}\|_2^2 = \sum_{i=1}^m w_i^2$  and  $\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_m \end{bmatrix} := \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}$

and  $z := b + \mathbf{x}^{(i)} \mathbf{w} = \mathbf{x}^{(i)} \boldsymbol{\theta}$



## $\ell_2$ - regularization

$$L(\boldsymbol{\theta}) = - \sum_{i=1}^n \left[ y^{(i)} \log(\sigma(z)) + (1 - y^{(i)}) \log(1 - \sigma(z)) \right] + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

- Now **two** objectives:      Minimize the “distance” to the data      but **also**      keep the **weights small**
- Other classifiers and regression models can be regularized in the same way!



## $\ell_2$ - regularization

### Regularisation parameter $\lambda$

- Provides **control** over **how well** the training data is fitted while keeping the weights small
- **Increasing** the value of  $\lambda \rightarrow$  **increases** the regularisation **strength**
- The `LogisticRegression` class in scikit-learn implements a parameter `C` for regularisation
- `C` is the inverse of  $\lambda$  (`C` is the inverse regularisation parameter):  $C := \frac{1}{\lambda}$
- Decreasing the value of `C`  $\rightarrow$  **increases** regularisation strength



## Train a logistic regression model

### Code example:

03\_logreg\_iris.ipynb

- Use the **iris** data set in scikit-learn
- Use ALL features
- Split the data into **training** and **test** set (`test_size=0.3, random_state=1`)
- Initialise the `LogisticRegression` class with `LogisticRegression(C=100.0, random_state=1)`
- Print out the number of **misclassified samples**
- Print out the **classification accuracy** for the **training** and **test data**

## $\ell_2$ - regularized logistic regression

$$+ \frac{\lambda}{2} ||\mathbf{w}||_2^2$$



### Code example:

- Use the **Wisconsin breast cancer** data set in scikit-learn
- Use ALL features
- Split the data 100 times into different training and test sets (test\_size=0.3)
- Initialise the `LogisticRegression` class with `LogisticRegression(C=100.0, random_state=1)`
- Print out the average and standard deviation of the validation accuracy across the 100 splits

03\_logreg\_cancer\_variance.ipynb



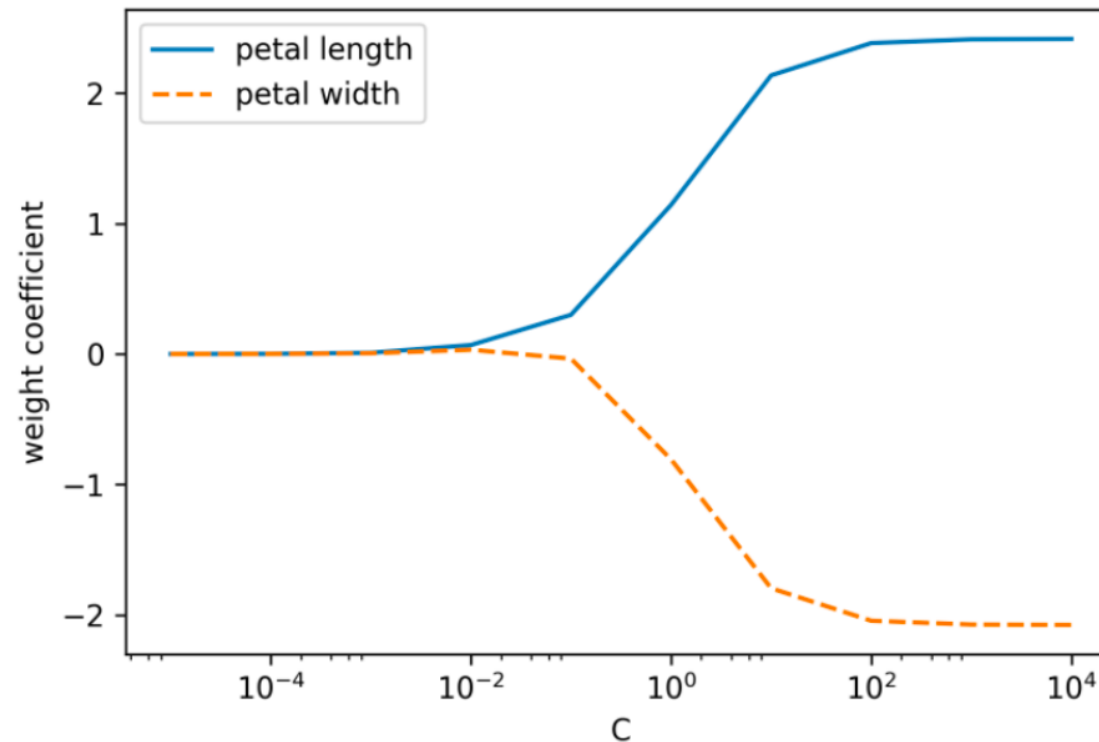
## Train a logistic regression model on the cancer data set

### Exercise:

- Use the **Wisconsin breast cancer** data set in scikit-learn
- Use ALL features
- Split the data into **training** and **test** set (`test_size=0.3, random_state=1`)
- Initialise the `LogisticRegression` class with `LogisticRegression(C=100.0, random_state=1)`
- Print out the number of **misclassified samples**
- Print out the **classification accuracy** for **training** and **test data**
- **Does the accuracy change with C?**

## $\ell_2$ - regularized logistic regression

The effect of regularisation can be visualised by plotting the “ $\ell_2$ -regularization path”



$$+ \frac{\lambda}{2} ||\mathbf{w}||_2^2$$



03\_logreg\_iris\_regularization.ipynb

Figure from Python Machine Learning (Raschka & Mirjalili)



## $\ell_2$ - regularized logistic regression

$$+ \frac{\lambda}{2} ||\mathbf{w}'||_2^2$$



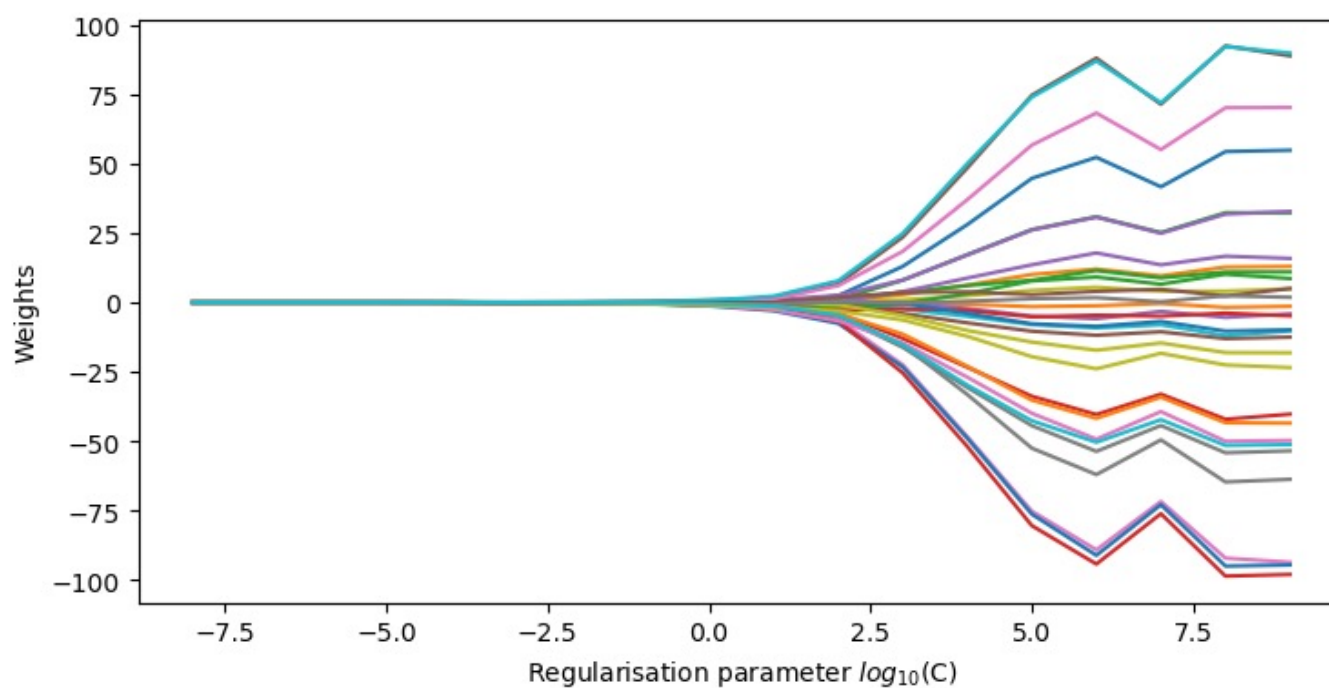
### Exercise:

- Use the **Wisconsin breast cancer** data set in scikit-learn
- Use ALL features
- Split the data into a training set and a test set (test\_size=0.3)
- Initialise the `LogisticRegression` class with `LogisticRegression(C=10.0**c, random_state=1)`
  - Small  $c$  varies from -8 to 10 (steps of 1)
- Plot the training and test accuracy across each  $C$  in one plot
- For which  $C$  does the model provide the best test accuracy? For which  $C$  values does the model overfit? For which  $C$  values does the model underfit?

03\_logreg\_cancer\_regularization.ipynb

# $\ell_2$ - regularized logistic regression

$$+ \frac{\lambda}{2} ||\mathbf{w}||_2^2$$

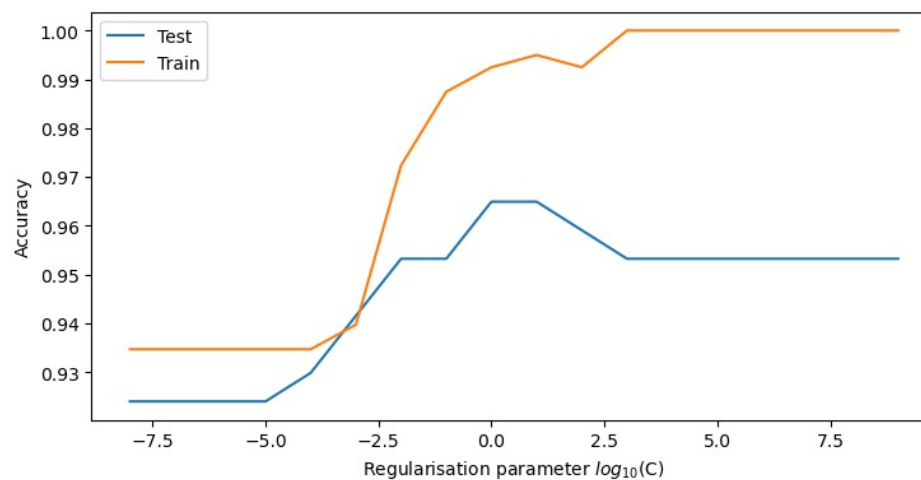
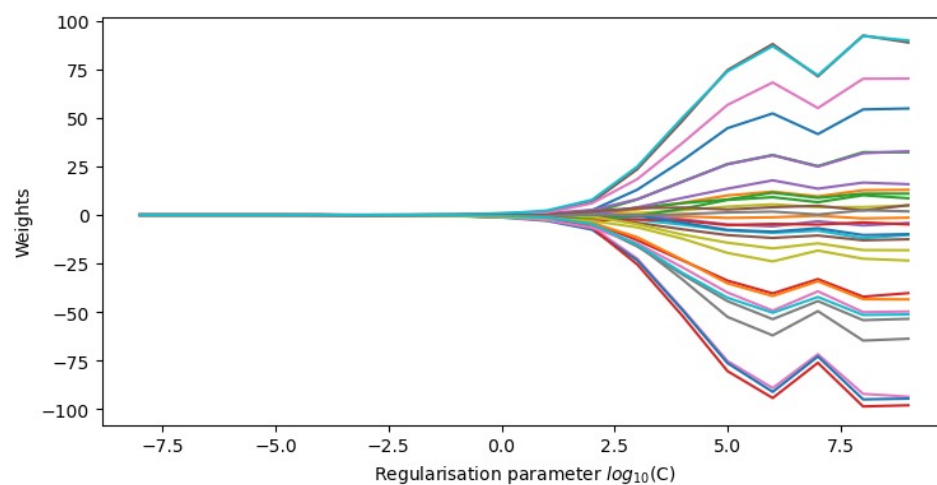


03\_logreg\_cancer\_regularization.ipynb

Figure from Python Machine Learning (Raschka & Mirjalili)



## $\ell_2$ -regularized logistic regression (exercise)



# Bias – Variance trade-off

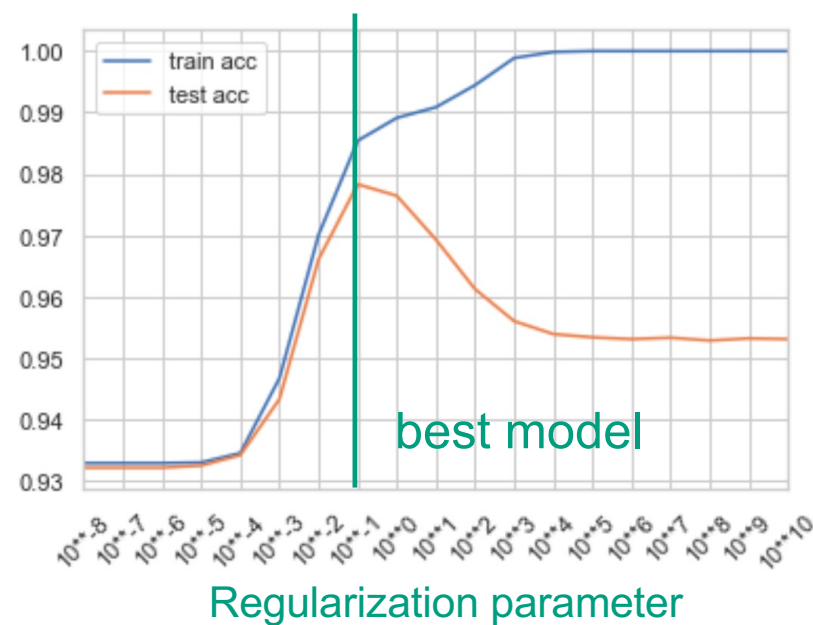
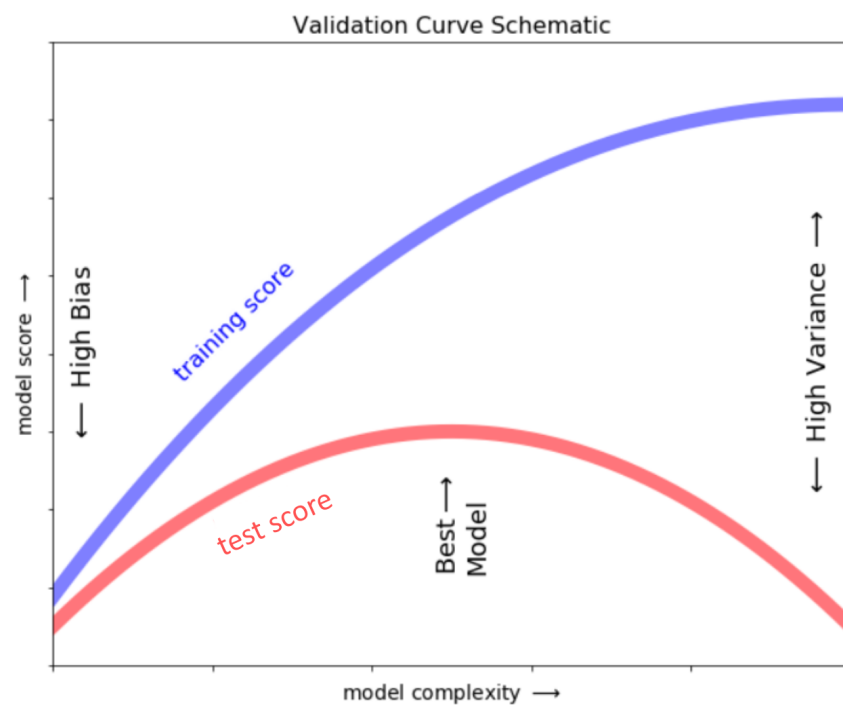
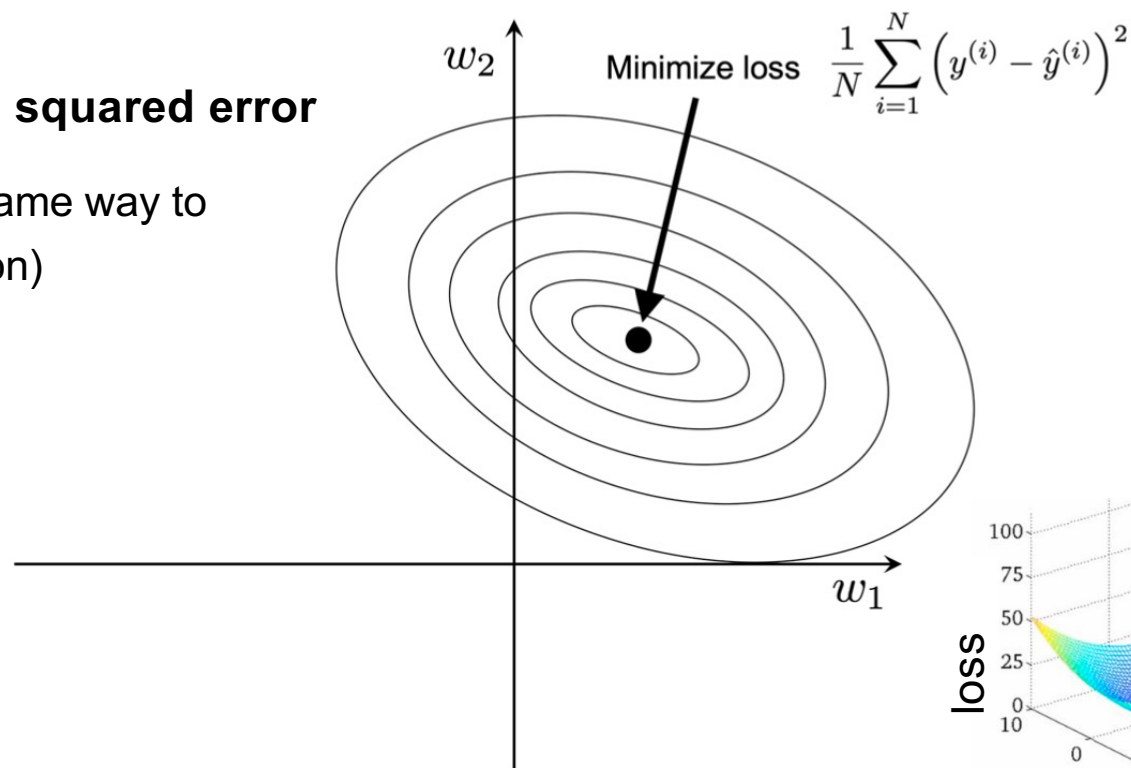


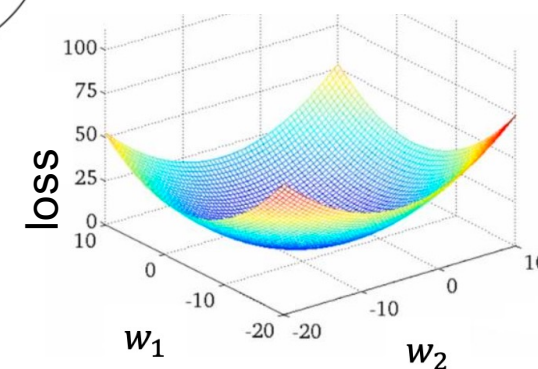
Figure from Python Machine Learning (Raschka & Mirjalili)

# Geometric interpretation of regularization

- Example: **Mean squared error**  
(Applies in the same way to logistic regression)
- Two features

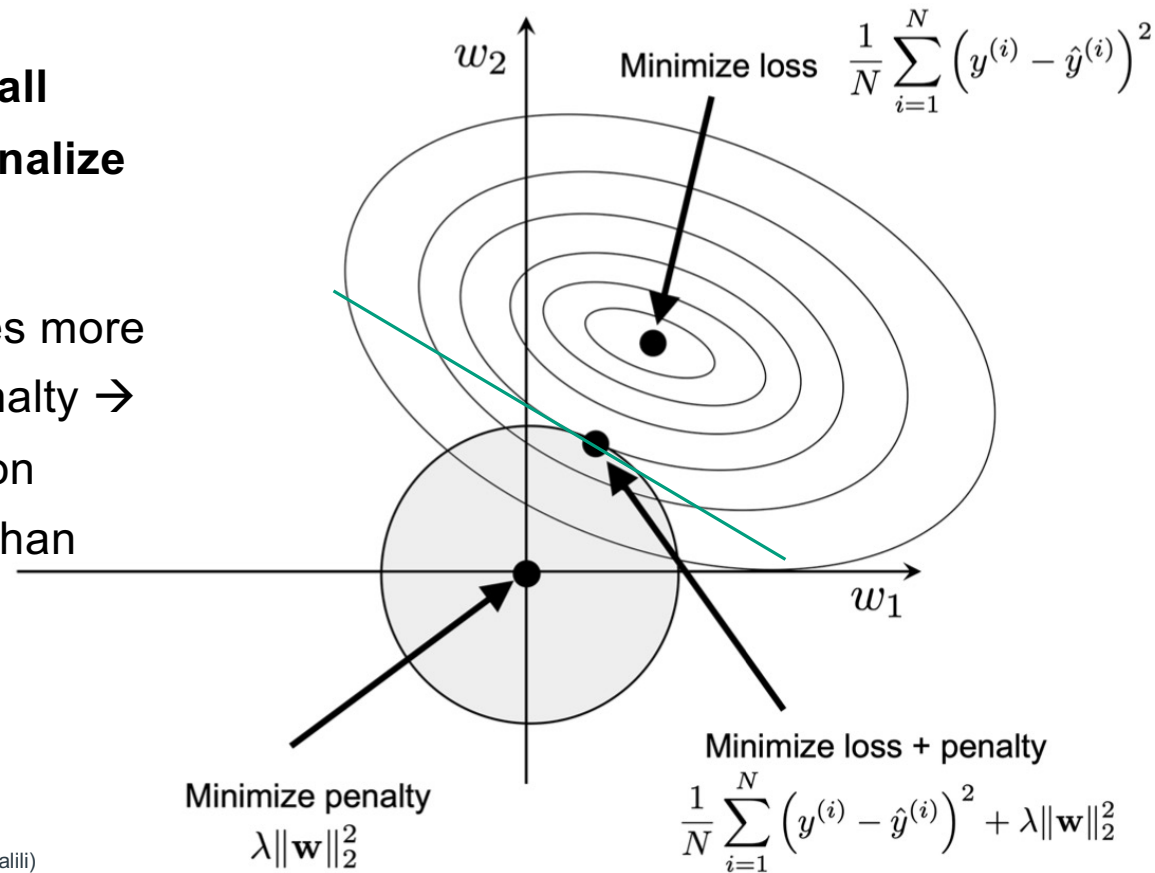


3D visualization



## Geometric interpretation: $\ell_2$ - regularization

- “Encourage small weights” or “Penalize large weights”
- Increasing  $\lambda$  gives more weight to the penalty  $\rightarrow$  More emphasis on smaller weights than fitting data





# Tackling overfitting via regularization

Sparsity-promoting regularization



$\ell_1$ - regularization (or L1 regularization / LASSO)

LASSO – Least absolute shrinkage and selection operator

$$L(\boldsymbol{\theta}) = - \sum_{i=1}^n \left[ y^{(i)} \log(\sigma(z)) + (1 - y^{(i)}) \log(1 - \sigma(z)) \right] + \lambda \|\mathbf{w}\|_1$$

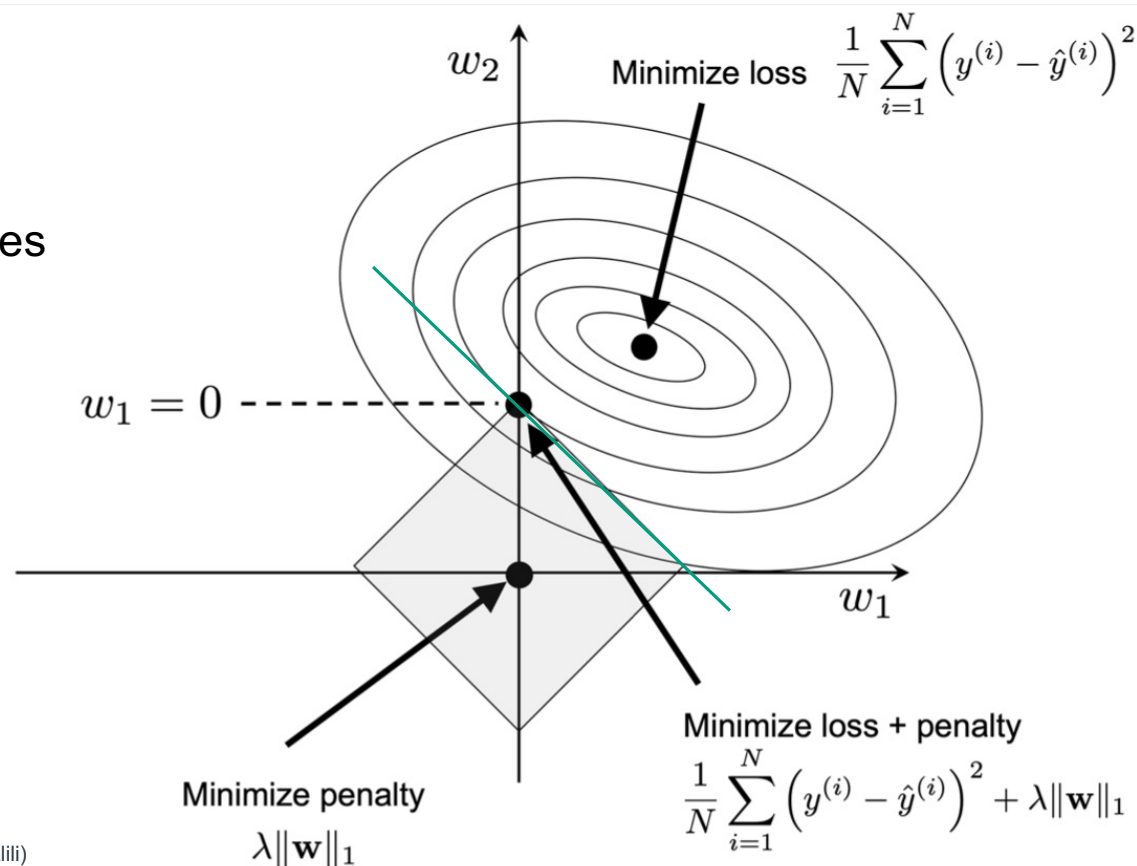
Where  $\|\mathbf{w}\|_1 = \sum_{i=1}^m |w_i|$  and  $\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_m \end{bmatrix} := \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}$

and  $z := b + \mathbf{x}^{(i)} \mathbf{w} = \mathbf{x}^{(i)} \boldsymbol{\theta}$



## Geometric interpretation: $\ell_1$ - regularization

- Example: Mean squared error (the same concept applies to the logistic loss function)

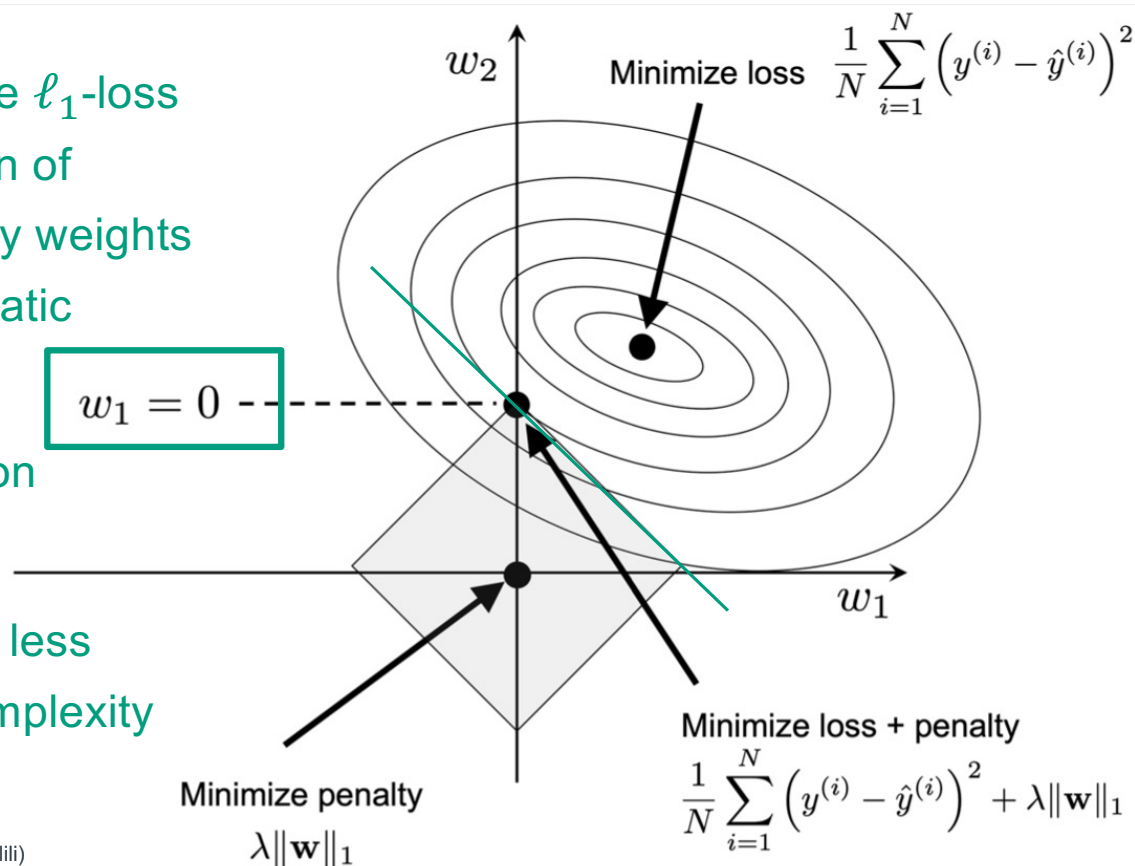


## Geometric interpretation: $\ell_1$ - regularization

Due to the shape of the  $\ell_1$ -loss this leads to a selection of weights such that many weights are zero (here: “automatic feature selection”)

We say  $\ell_1$  regularization promotes “**sparsity**”

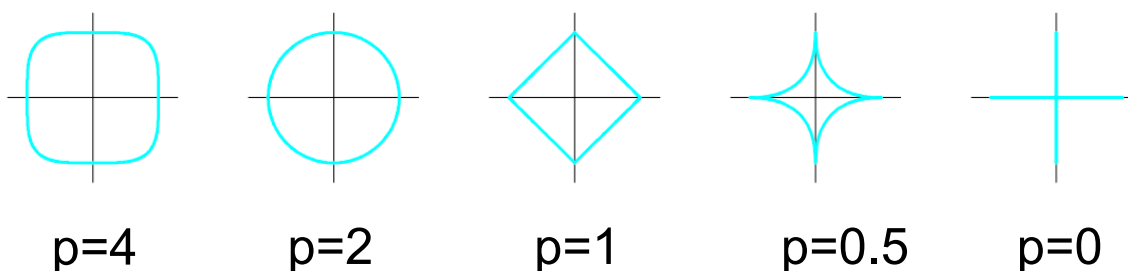
It leads to models with less parameters / lower complexity





## Geometric interpretation: $\ell_p$ - regularization

- Other theoretical options:



$$+ \frac{\lambda}{p} \|\mathbf{w}\|_p^p$$

$$\|\mathbf{w}\|_p^p = \sum_{i=1}^m |w_i|^p$$

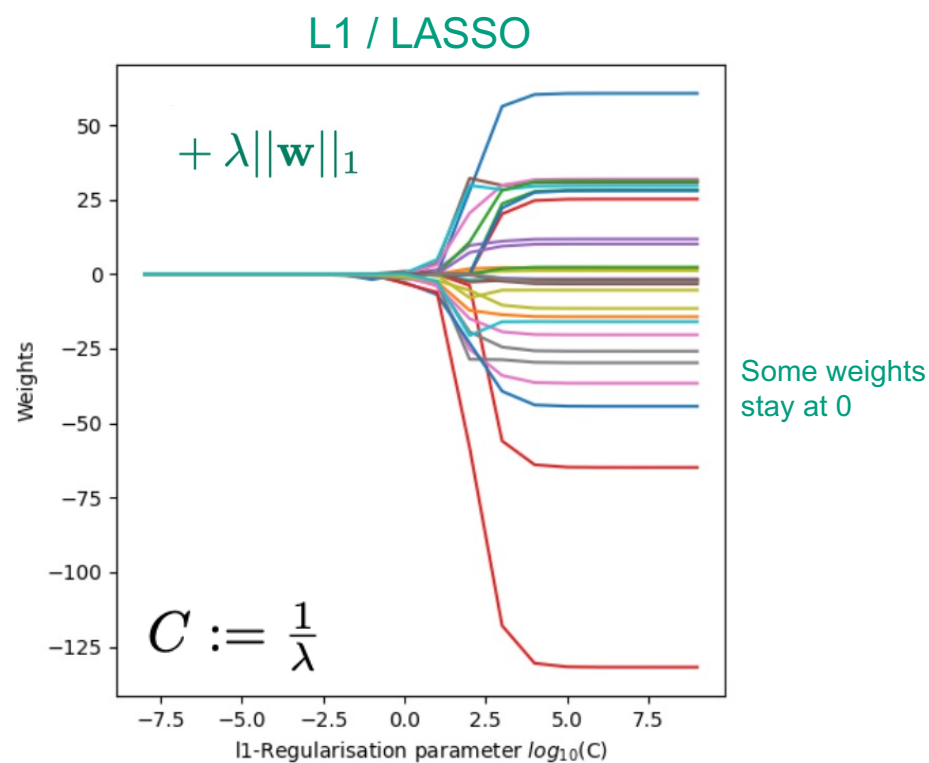
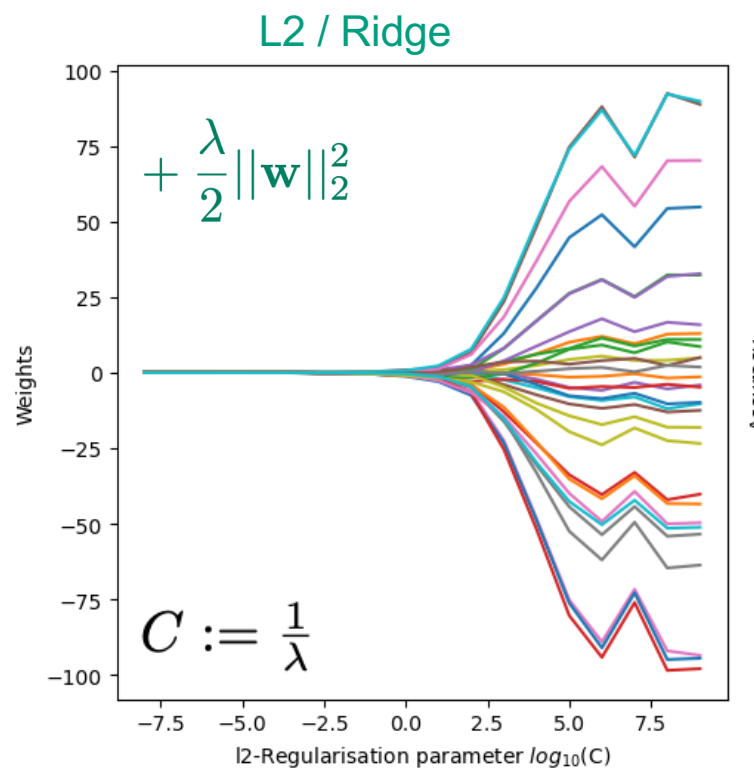
- Due to computational restrictions: **In practice, the options used are  $\ell_1$ (Lasso),  $\ell_2$ (Ridge) and  $\ell_1 + \ell_2$  (Elastic net)**



## Sparsity-promoting regularization

- L1 regularisation usually yields **sparse** feature vectors → most feature weights will be zero
- Sparsity can be useful in practice if we have a **high-dimensional** dataset with **many** features that are **irrelevant**
- L1 regularisation is especially useful in cases where more **irrelevant dimensions** than **samples** are present
- L1 regularisation can be understood as a technique for **feature selection**

# Comparison of L1/L2 regularization



`03_logreg_cancer_regularization_L1_L2.ipynb`



## Scikit-learn classifiers with a regularization option

- Perceptron
- Logistic regression
- Linear SVC

### Options:

- L1 – a.k.a LASSO (Least Absolute shrinkage and Selection Operator)
- L2 – a.k.a. Ridge penalisation
- L1 + L2 – a.k.a. Elastic Net

