

Data Science and Machine Learning in Python

Stephan Weyers

Part 1: Data Science

	Date	Topics covered
1	Apr 13 th	Course introduction Data Science motivation How to use Jupyter Notebook Python types and lists Loops, if/else, functions
2	Apr 20 th	Python tuples, lists, dictionaries Functions Numpy basics, operations Image processing
3	Apr 27 th	Pandas Series, DataFrame Pandas basic operations Import/export files
4	May 4 th	Principles of data visualization Data cleaning and preparation Join, combine and reshape data
5	May 11 th	Volkswahl Bund dataset Data visualization in Python How to write Data Science reports Data aggregation and grouping

Part 2: Machine Learning

	Date	Topics covered
6	Jun 1 st	Introduction to supervised learning Classification and regression scikit-learn k-Nearest Neighbors Linear regression (ridge and lasso)
7	Jun 8 th	Linear classification models Decision trees Random forests and gradient boosting
8	Jun 15 th	Kernel support vector machines Neural networks
9	Jun 22 nd	Introduction to unsupervised learning Preprocessing and scaling Dimensionality reduction Principal component analysis
10	Jun 29 th	k-means clustering Hierarchical clustering DBSCAN
11	Jul 6 th	Representing data Engineering features Model evaluation and improvement Text data analysis

Deadlines for Submission and Distribution of Grading

Student task	Deliverables	Deadline	Work	Share of grade
W01 Assignment	Code and results	Apr 26 th	Team A	5.0%
W02 Case Study	Code / presentation slides	May 22 nd	Team B	18.0%
W02 Case Study	Peer review*	May 31 st	Individual	2.0%
W03 Assignment	Code and results	May 29 th	Team B	5.0%
W04 Assignment	Code and results	Jun 12 th	Team C	10.0%
W05 Assignment	Code and results	Jun 26 th	Team D	7.0%
W06 Assignment	Code and results	Jul 8 th	Team D	13.0%
W07 Case Study	Code / presentation slides	Jul 17 th	Team D	22.0%
W07 Case Study	Peer review*	Jul 31 st	Individual	3.0%
DataCamp 1	Finish course	May 9 th	Individual	2.5%
DataCamp 2	Finish course	May 30 th	Individual	2.5%
DataCamp 3	Finish course	Jun 20 th	Individual	2.5%
DataCamp 4	Finish course	Jul 11 th	Individual	2.5%

* Peer review is mandatory. Quality of peer review itself is graded. Not providing peer review at all would result in high point deduction

Agenda for online lecture 8

Session	Topic	Mode	Materials used	Minutes	End
14:30-16:00	Organizational questions	Q&A		10	14:40
	Bank Marketing data	Team work in break-out rooms	Lecture 08a notebook	45	15:25
	Kernel SVM	Lecture / Q&A	Lecture slides	15	15:40
	Neural networks	Lecture / Q&A	Lecture slides	15	15:55
16:10-17:40	Telco Customer Churn	Lecture / Q&A	Lecture 08b notebook	30	16:40
	Online retail exercise 1+2	Team work in break-out rooms	Lecture 08c notebook	40	17:20
	Online retail exercise 1	Lecture / Q&A	Lecture 08c notebook	15	17:35
17:50-19:20	Online retail exercise 2+3	Team work in break-out rooms	Lecture 08c notebook	30	18:20
	Online retail exercise 2	Lecture / Q&A	Lecture 08c notebook	10	18:30
	Online retail exercise 3	Team work in break-out rooms	Lecture 08c notebook	30	19:00
	Online retail exercise 3	Lecture / Q&A	Lecture 08c notebook	10	19:10
	Organizational questions	Q&A		10	19:20

Supervised Approaches

- Labeled data
- Target values known

Classification

- Predict category

Regression

- Predict numeric value

Unsupervised Approaches

- Unlabeled data
- No target value provided

Cluster Analysis

- Organize similar cases into segments

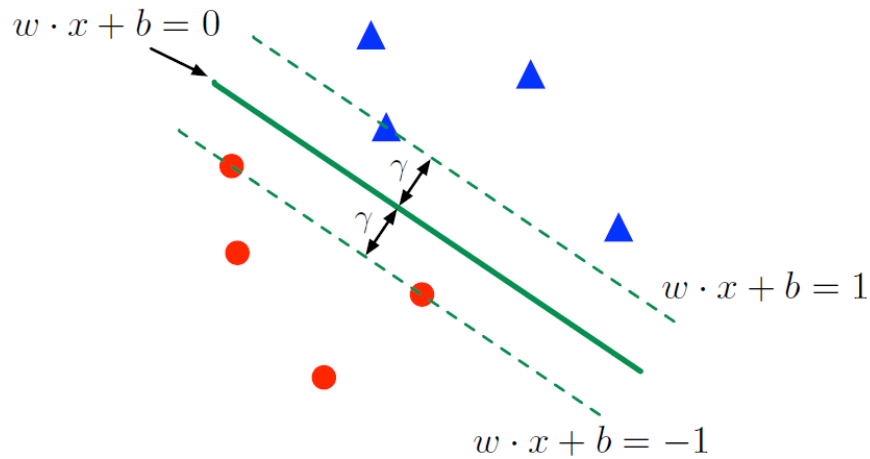
Dimensionality reduction

- Reduce number of features

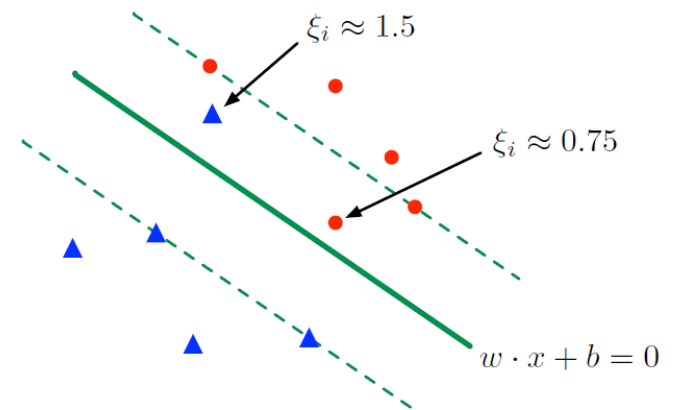
Question for discussion

- Find examples for each of the 4 categories

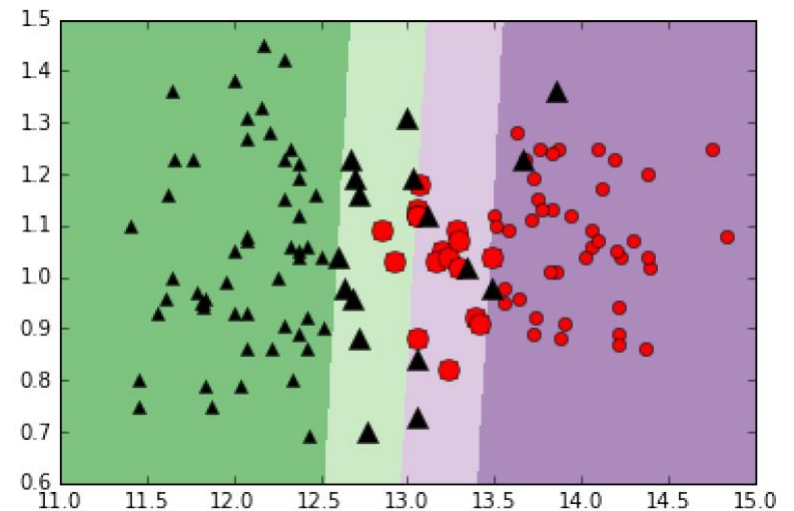
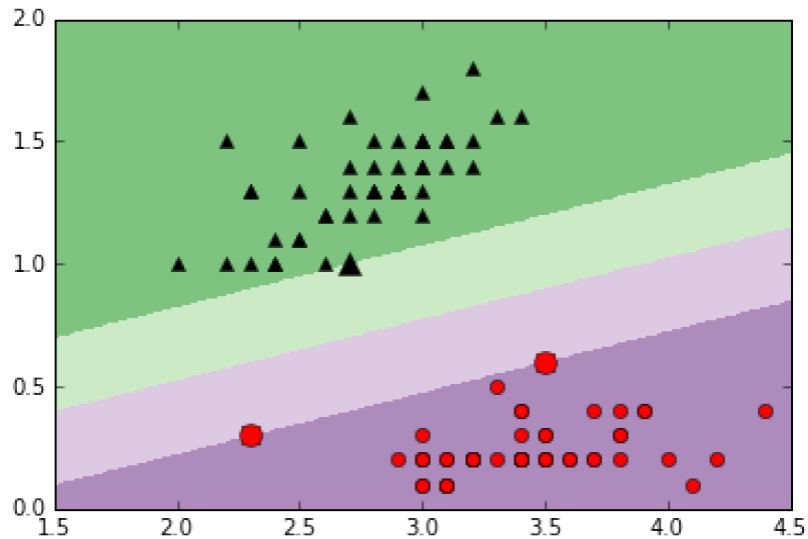
Maximize margin $\gamma = \frac{1}{\|w\|}$



Minimize slack ξ



Support vectors



Given training input data $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)} \in \mathbb{R}^d$ with labels $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(n)} \in \{-1, 1\}$.

Try to find $\mathbf{w} \in \mathbb{R}^d$ and $b \in \mathbb{R}$, so that

$$\text{sign}(\mathbf{w} \cdot \mathbf{x}^{(i)} + b) = \mathbf{y}^{(i)}$$

as often as possible.

Logistic regression

Minimize

$$L(\mathbf{w}, b) = -\ln \left(\prod_{i=1}^n \text{Pr}_{\mathbf{w}, b}(\mathbf{y}^{(i)} \mid \mathbf{x}^{(i)}) \right) + \lambda \|\mathbf{w}\|_2^2 = -\sum_{i=1}^n \ln \left(\frac{1}{1 + e^{-\mathbf{y}^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + b)}} \right) + \lambda \|\mathbf{w}\|_2^2$$

Linear support vector machines

Hard margin

$$\min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \|\mathbf{w}\|_2^2$$

such that

$$\mathbf{y}^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + b) \geq 1$$

Soft margin

$$\min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}, \xi \in \mathbb{R}^n} \|\mathbf{w}\|_2^2 + K \sum_{i=1}^n \xi_i$$

such that

$$\mathbf{y}^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)} + b) \geq 1 - \xi_i$$

$$\xi \geq 0$$

Parameters

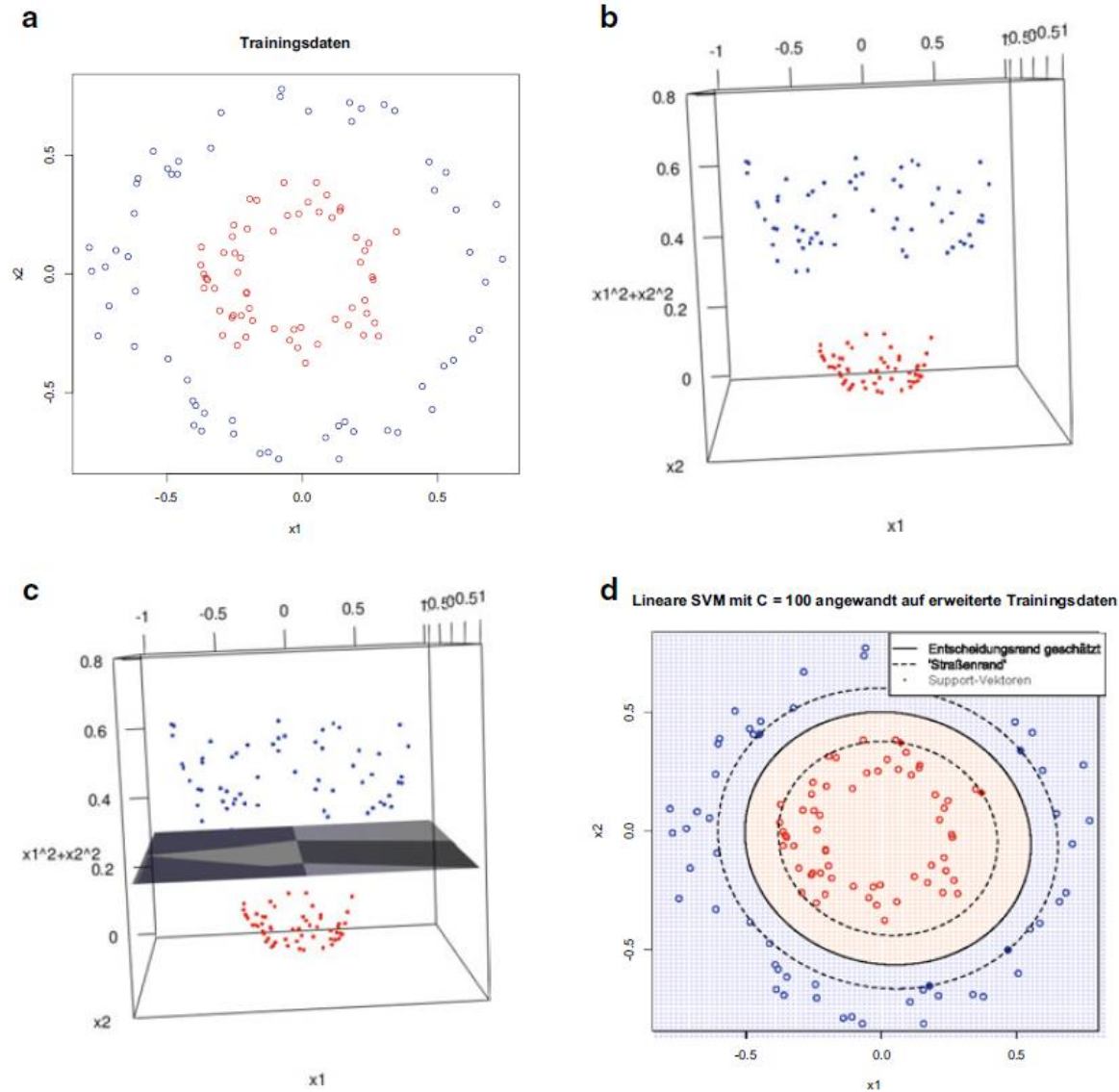
- Regularization parameter α and C
- Model type lasso vs. ridge for regression / logistic vs. SVM for classification

Strengths

- Fast to train, fast to predict
- Work well with sparse data
- Relatively easy to understand how predictions are made
- Work well with large number of features

Weaknesses

- Coefficients hard to interpret, especially if features are highly correlated
- Sometimes fail with small datasets
- Perform bad with non-linear features and datasets that are not linearly separable



Input variables

$$x = (x_1, x_2, x_3)$$

Extended input variables (degree p=2)

$$\phi(x) = (x_1, x_2, x_3, x_1^2, x_2^2, x_3^2, x_1x_2, x_1x_3, x_2x_3)$$

Extended input variables (degree p=3)

$$\phi(x) = (x_1, x_2, x_3, x_1^2, x_2^2, x_3^2, x_1x_2, x_1x_3, x_2x_3, x_1^3, x_2^3, x_3^3, x_1^2x_2, x_1^2x_3, x_2^2x_1, x_2^2x_3, x_3^2x_1, x_3^2x_2, x_1x_2x_3)$$

Kernel trick

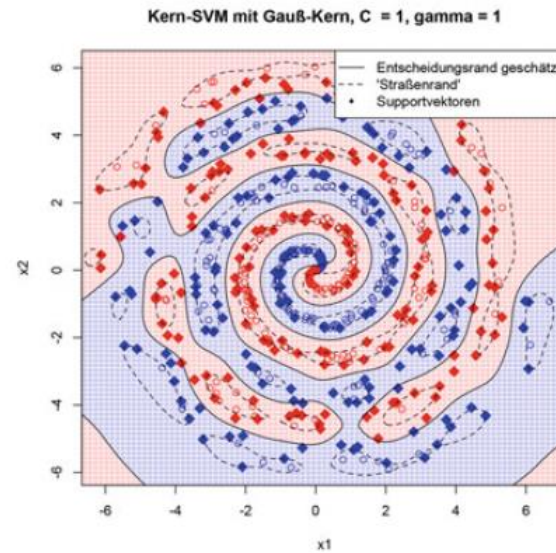
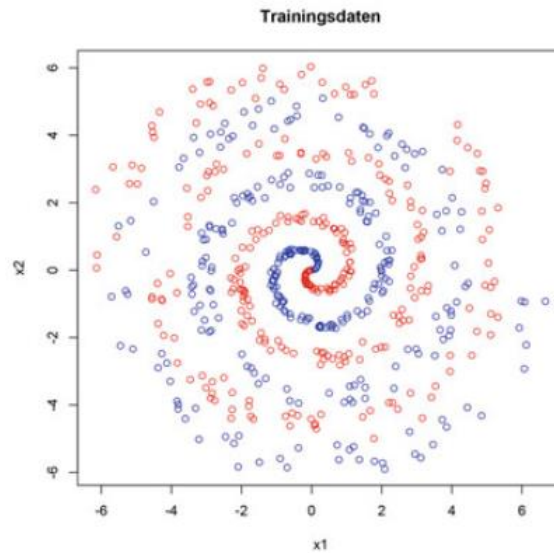
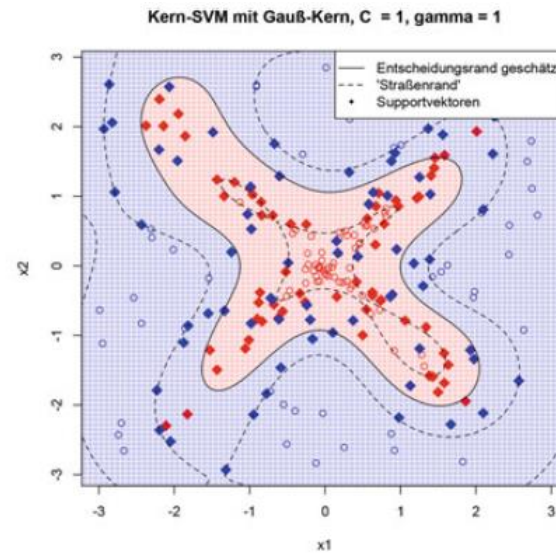
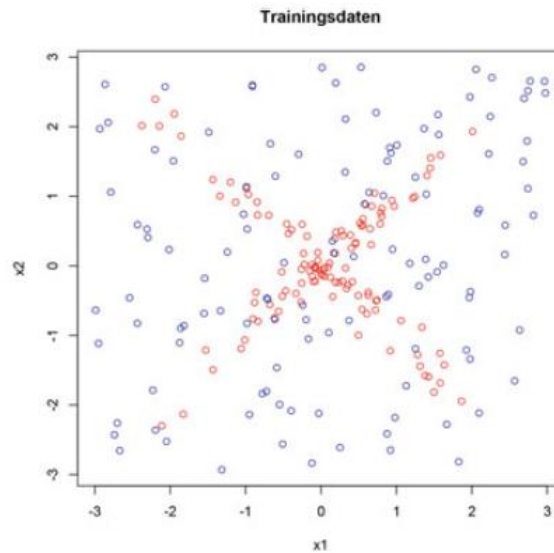
$$\phi(x) \cdot \phi(y) = (1 + x \cdot y)^p$$

Polynomial kernel

$$\phi(x) \cdot \phi(y) = (1 + x \cdot y)^p$$

Gaussian (RBF) kernel

$$\phi(x) \cdot \phi(y) = e^{-\gamma \|x-y\|^2}$$



Parameters

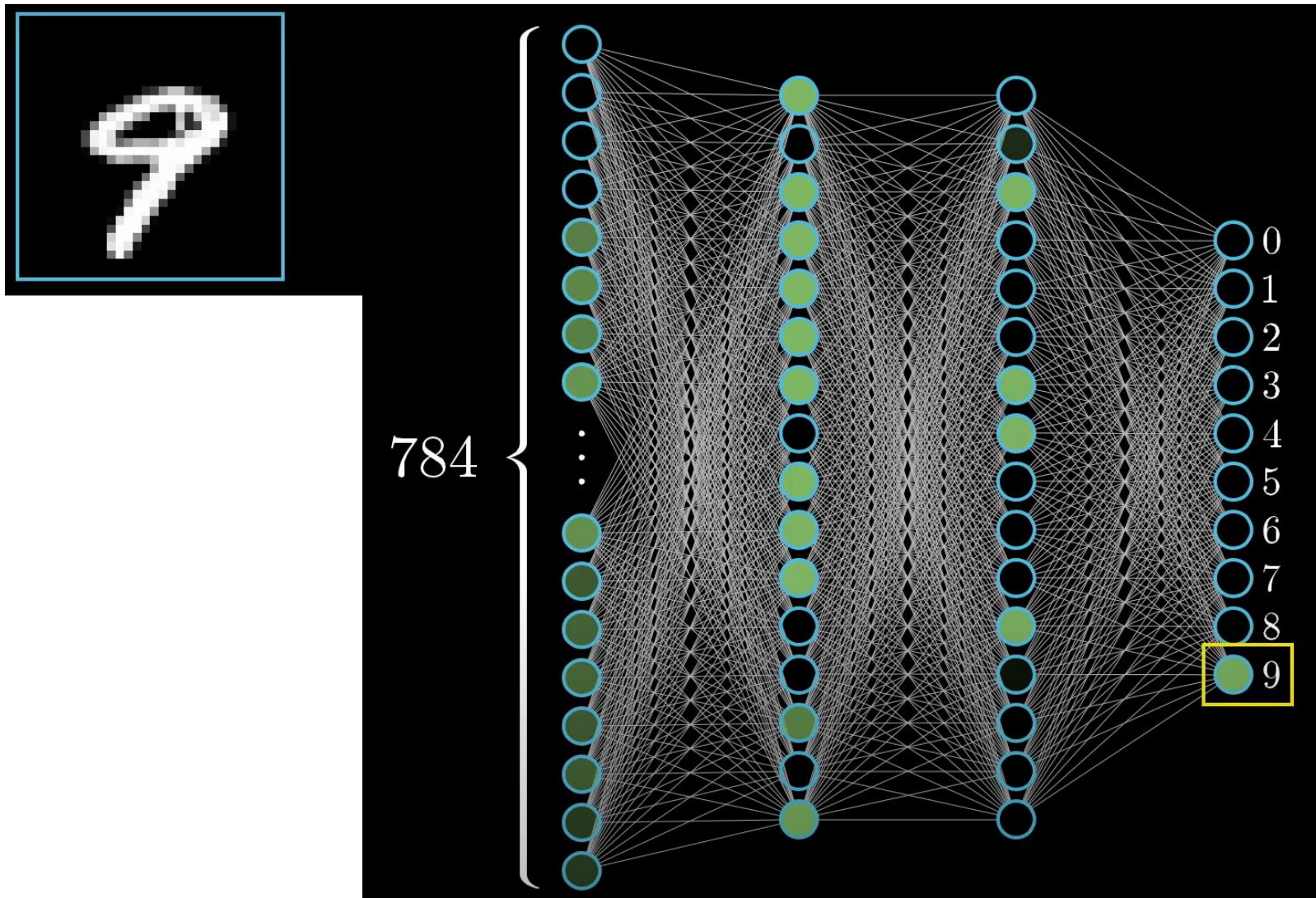
- Regularization parameter C
- Kernel “poly” or “rbf”
- Kernel specific parameters “degree” or “gamma”

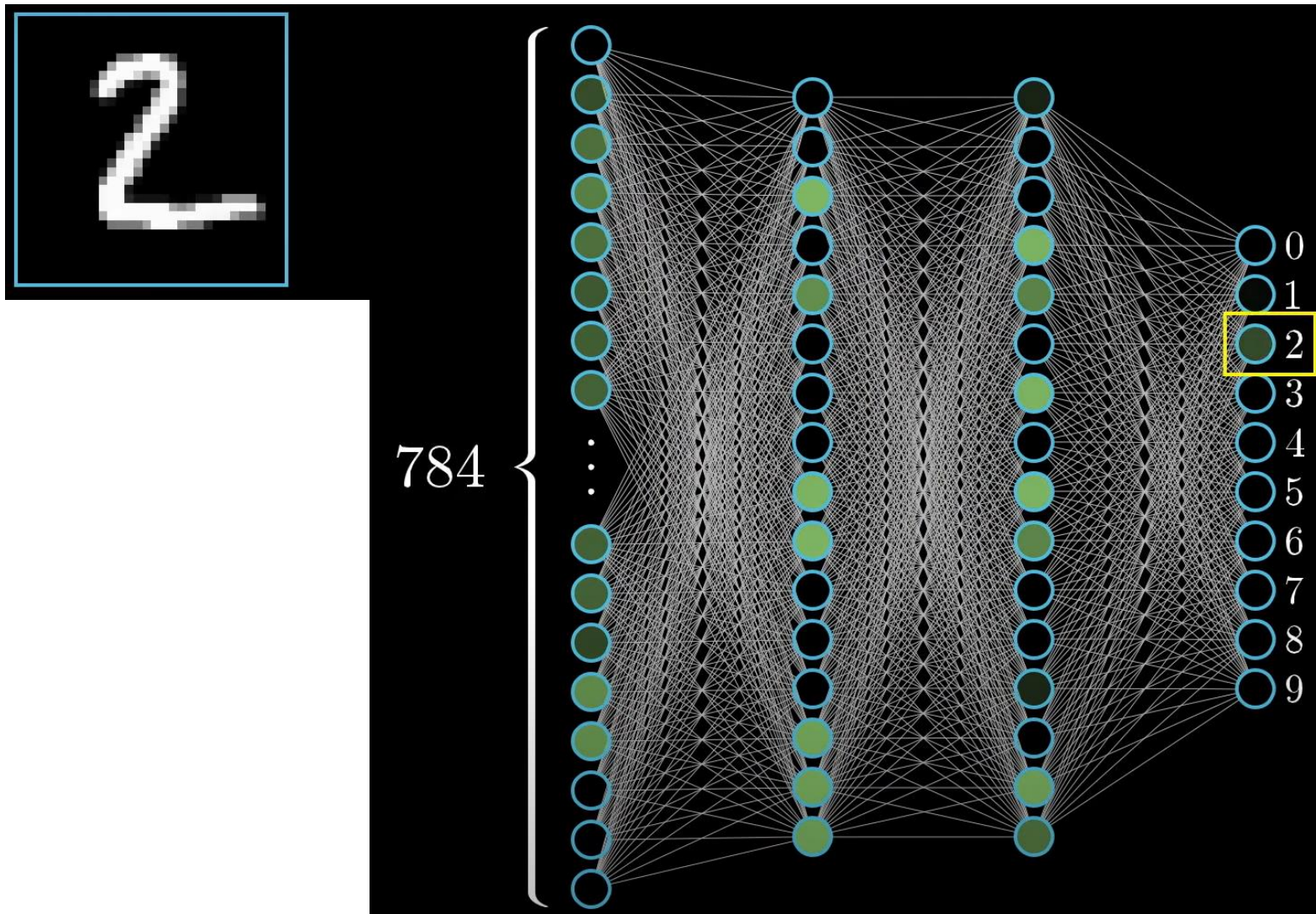
Strengths

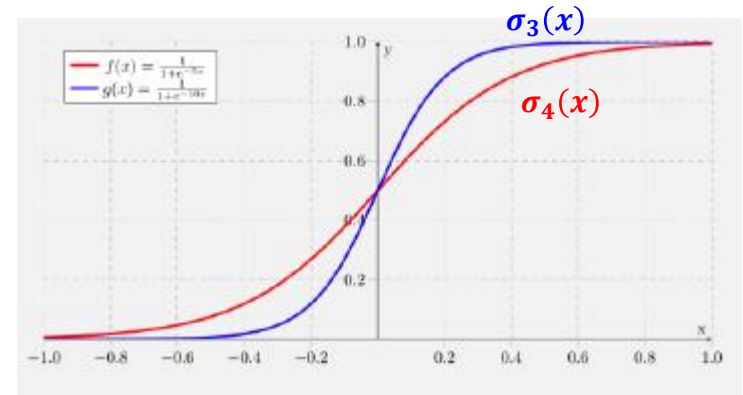
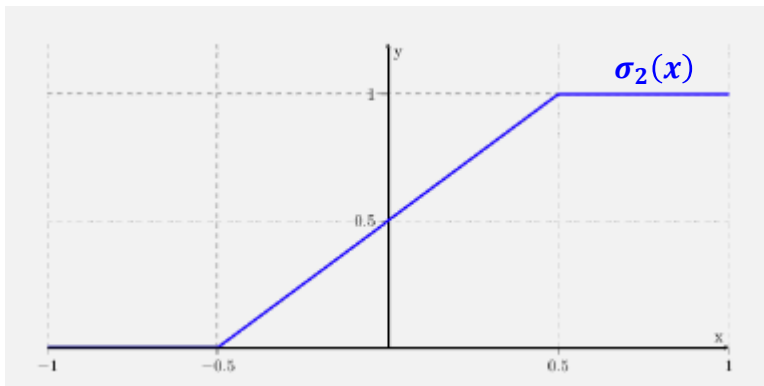
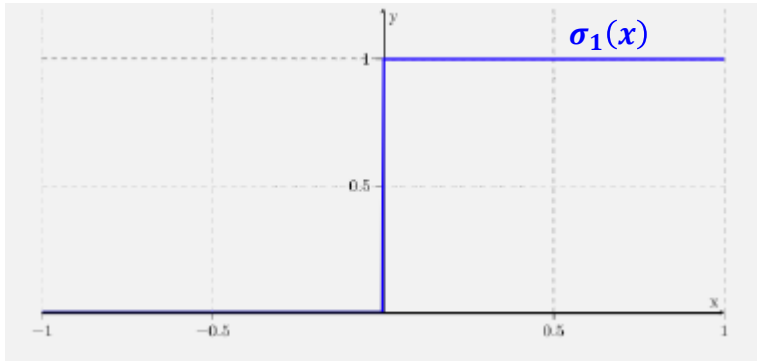
- Complex decision boundaries
- Work well with few and also many features
- Work well, if all features are on same scale / in same units

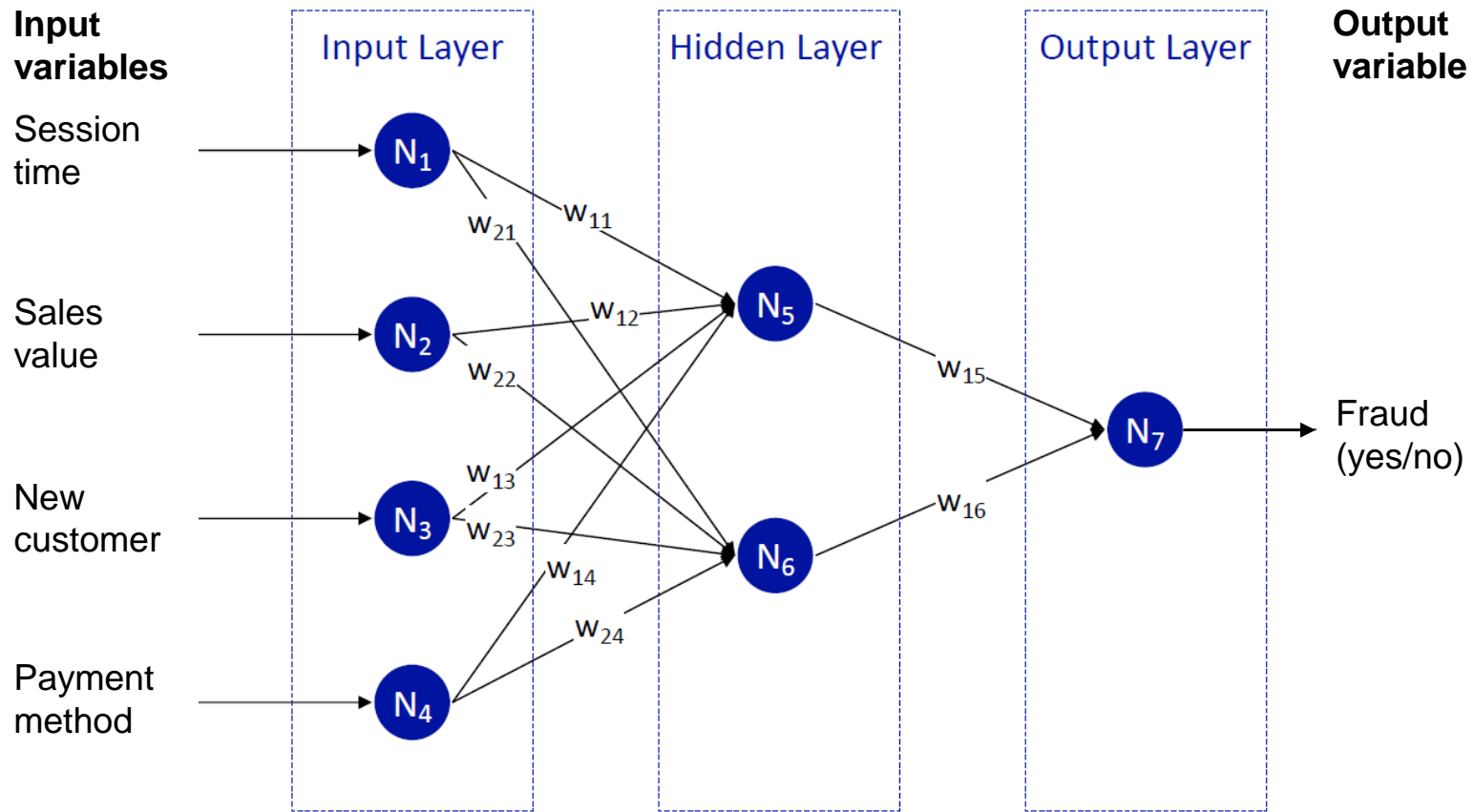
Weaknesses

- Runtime and memory usage challenges with high number of samples
- Require preprocessing / scaling of data to get good results
- Results hard to interpret, difficult to explain to non-experts





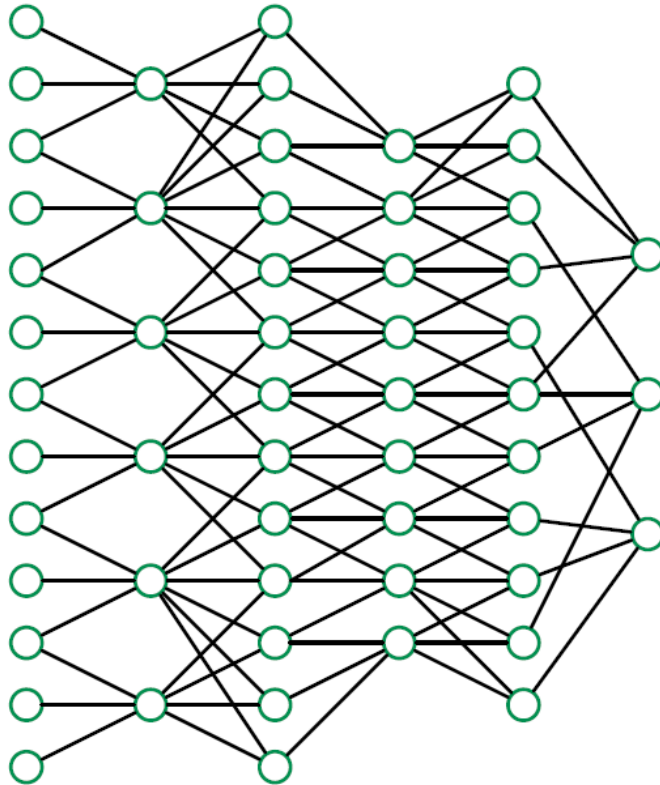




$$N_5 = \sigma(w_{11}N_1 + w_{12}N_2 + w_{13}N_3 + w_{14}N_4 + b_5)$$

$$N_6 = \sigma(w_{21}N_1 + w_{22}N_2 + w_{23}N_3 + w_{24}N_4 + b_6)$$

$$N_7 = \sigma(w_{15}N_5 + w_{16}N_6 + b_7)$$



Question for discussion

A neural network has 6 layers, each consisting of 100 nodes, and each fully connected to the previous layer. Roughly how many parameters (w_{ij}) does this network have?

Parameters

- Activation function, number and size of layers
- Many more

Strengths

- State-of-the-art in many applications
- Able to build very complex models

Weaknesses

- Black box – difficult to interpret and explain results
- Careful preprocessing and scaling of data needed
- Very long computation time for large datasets
- Difficulties with mixed kinds of features on different scales