

# **Topological Machine Learning: The (W)Hole Truth**

**Lecture 4**

Bastian Rieck (@Pseudomanifold)

# Preliminaries

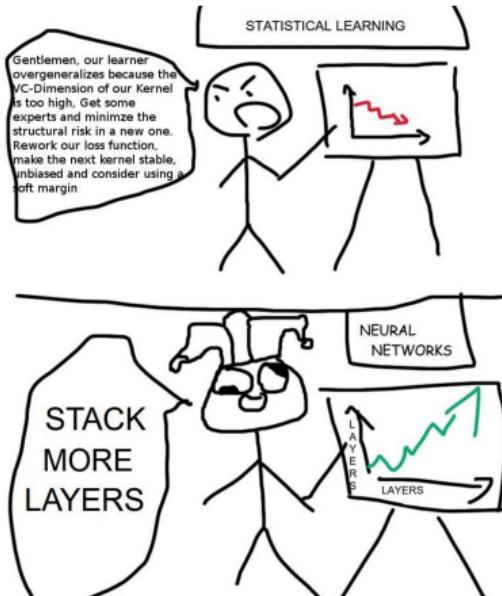
Do you have feedback or any questions? Write to [bastian.rieck@helmholtz-muenchen.de](mailto:bastian.rieck@helmholtz-muenchen.de) or reach out to @Pseudomanifold on Twitter. You can find the slides and additional information with links to more literature here:

<https://heidelberg.topology.rocks>

# What is machine learning?



# Two households, both alike in dignity...



# Two households, both alike in dignity...

## Classical machine learning

- ☆ Somewhat better interpretable
- ☆ Works well in small-sample regime
- ☆ Low requirements on compute

## Deep learning

- ☆ Often like a ‘black box’
- ☆ Requires large amounts of data
- ☆ High compute requirements

Classical models are often not easy to train because they require ‘hand-crafted’ features. Deep models, on the other hand, are harder to ‘investigate’ *a posteriori*.

# When to pick one over the other?

Two anecdotes

In particular we ask whether supervised deep learning can be used for cell annotation, i.e. to predict cell-type labels from single-cell gene expression profiles. *After evaluating 10 classification methods across 14 datasets, we notably find that deep learning does not outperform classical machine-learning methods in the task.*

(N. D. Köhler, M. Büttner, N. Andriamanga and F. J. Theis, 'Deep learning does not outperform classical machine learning for cell-type annotation', Preprint, 2021)

*One of the most surprising results is the highly-competitive performance of our DTW-KNN ensemble classifier, whose performance for earlier horizons exceeds all of the other methods, despite its conceptual simplicity.* While this is highly relevant for the early detection of sepsis, the DTW-KNN classifier suffers from some practical limitations that impede online monitoring scenarios and its application to very large patient cohorts.

(M. Moor, M. Horn, **B. Rieck**, D. Roqueiro and K. Borgwardt, 'Early Recognition of Sepsis with Gaussian Process Temporal Convolutional Networks and Dynamic Time Warping', *Proceedings of the 4th Machine Learning for Healthcare Conference*, ed. by F. Doshi-Velez et al., Proceedings of Machine Learning Research 106, PMLR, 2019, pp. 2–26, arXiv: 1902 . 01659 [cs.LG])

A brief overview of some tried and tested classical  
machine learning methods

# Starting with the data

## 'Wine' data set

- ☆ 178 samples
- ☆ 13 attributes
- ☆ 3 classes, corresponding to cultivars

# Starting with the data

## 'Wine' data set

- ☆ 178 samples
- ☆ 13 attributes
- ☆ 3 classes, corresponding to cultivars

Normally, we should try to understand all of the attributes: are they *continuous* or *discrete*? What are their ranges?

# Starting with the data

## 'Wine' data set

- ☆ 178 samples
- ☆ 13 attributes
- ☆ 3 classes, corresponding to cultivars

Normally, we should try to understand all of the attributes: are they *continuous* or *discrete*? What are their ranges?

## Attributes

- ☆ Alcohol
- ☆ Malic acid
- ☆ Ash
- ☆ Alcalinity of ash
- ☆ Magnesium
- ☆ Total phenols
- ☆ Flavanoids
- ☆ Nonflavanoid phenols
- ☆ Proanthocyanins
- ☆ Colour intensity
- ☆ Hue
- ☆ OD<sub>280</sub>/OD<sub>315</sub> value of diluted wines
- ☆ Proline

# Starting with the data

## 'Wine' data set

- ☆ 178 samples
- ☆ 13 attributes
- ☆ 3 classes, corresponding to cultivars

Normally, we should try to understand all of the attributes: are they *continuous* or *discrete*? What are their ranges?

## Attributes

- ☆ Alcohol
- ☆ Malic acid
- ☆ Ash
- ☆ Alcalinity of ash
- ☆ Magnesium
- ☆ Total phenols
- ☆ Flavanoids
- ☆ Nonflavanoid phenols
- ☆ Proanthocyanins
- ☆ Colour intensity
- ☆ Hue
- ☆ OD280/OD315 value of diluted wines
- ☆ Proline

Source: <https://archive.ics.uci.edu/ml/datasets/wine>

# Logistic regression

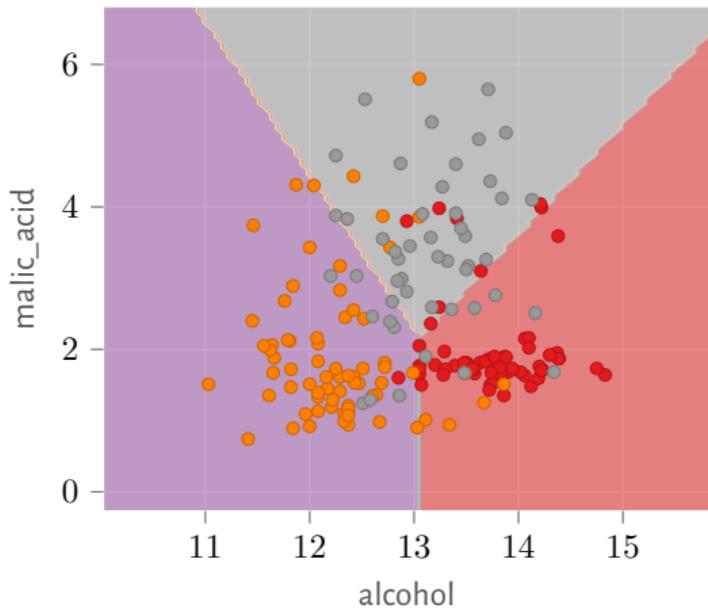
## Brief description

Logistic regression models the probability of an event taking place by having the log-odds for the event be a *linear combination* of one or more independent variables. We then estimate the parameters, i.e. the coefficients, of such a model.

# Logistic regression

## Brief description

Logistic regression models the probability of an event taking place by having the log-odds for the event be a *linear combination* of one or more independent variables. We then estimate the parameters, i.e. the coefficients, of such a model.

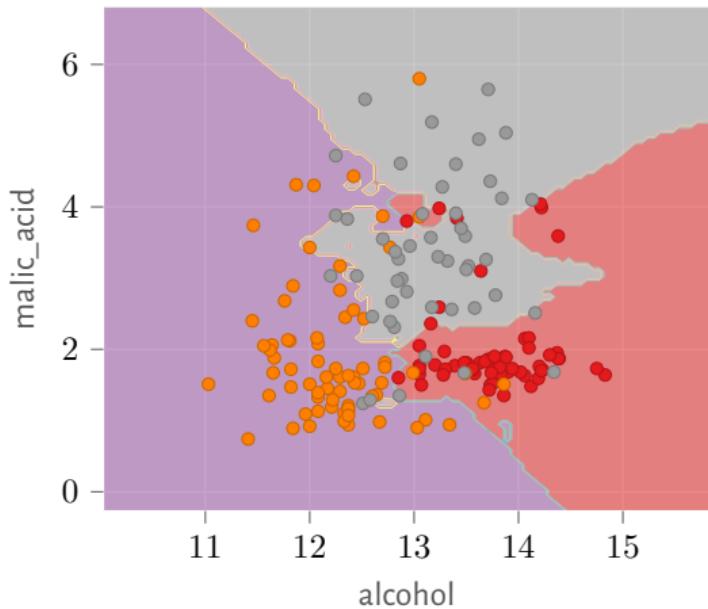


# $k$ -nearest neighbours classification and regression

- ☆ The quintessential ‘lazy learner’
- ☆ Requires only a similarity measure and a neighbourhood choice
- ☆ Surprisingly effective and efficient

# $k$ -nearest neighbours classification and regression

- ☆ The quintessential ‘lazy learner’
- ☆ Requires only a similarity measure and a neighbourhood choice
- ☆ Surprisingly effective and efficient



# Support vector machines

Given data points  $\mathbf{x}_i$  with labels  $y_i \in \{+1, -1\}$ ,  
find weights  $\mathbf{w}$  that minimise  $\lambda \|\mathbf{w}\|^2 +$   
 $[\frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i - b))]$ .

# Support vector machines

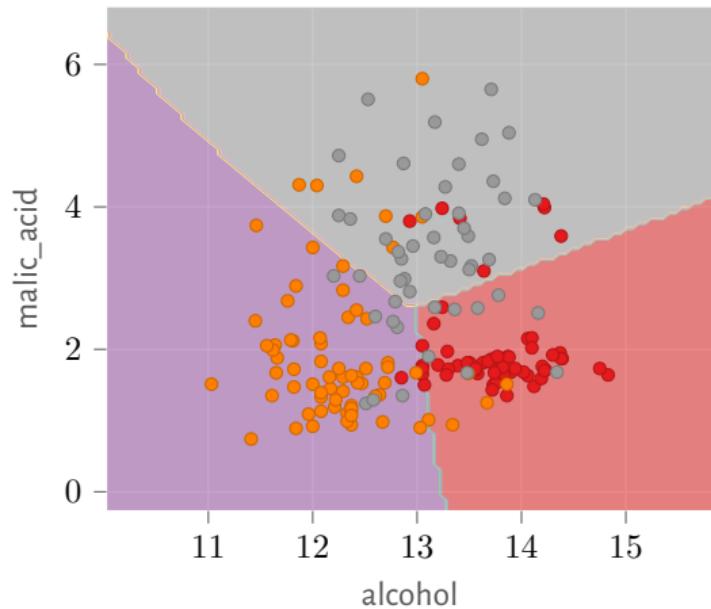
Given data points  $\mathbf{x}_i$  with labels  $y_i \in \{+1, -1\}$ ,  
find weights  $\mathbf{w}$  that minimise  $\lambda \|\mathbf{w}\|^2 +$   
 $[\frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i - b))]$ .

SVMs also support *kernels*, i.e. domain-specific  
similarity measures!

# Support vector machines

Given data points  $\mathbf{x}_i$  with labels  $y_i \in \{+1, -1\}$ ,  
find weights  $\mathbf{w}$  that minimise  $\lambda \|\mathbf{w}\|^2 + [\frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i - b))]$ .

SVMs also support *kernels*, i.e. domain-specific similarity measures!



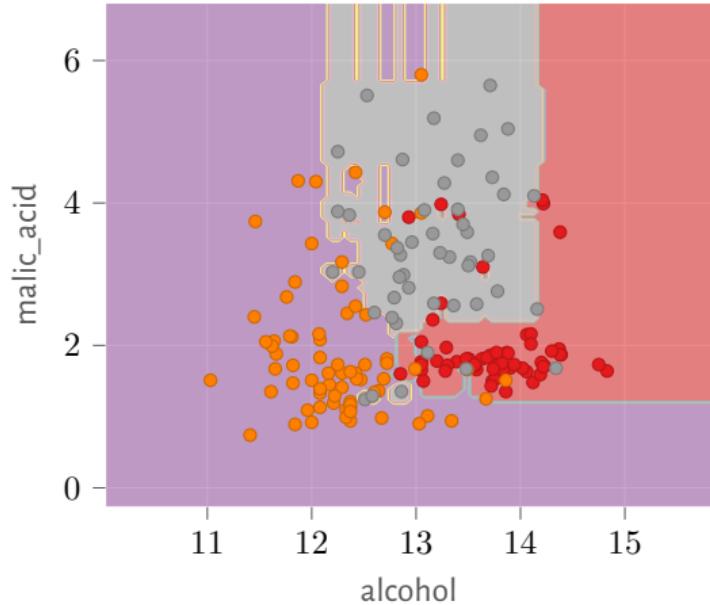
# Random forests

- ☆ Learn many small *decision trees* and use their joint predictions
- ☆ Leads to *highly interpretable* models!
- ☆ Still an active area of research<sup>1</sup>

<sup>1</sup>G. Ke et al., ‘LightGBM: A Highly Efficient Gradient Boosting Decision Tree’, *Advances in Neural Information Processing Systems*, ed. by I. Guyon et al., vol. 30, Curran Associates, Inc., 2017

# Random forests

- ☆ Learn many small *decision trees* and use their joint predictions
- ☆ Leads to *highly interpretable* models!
- ☆ Still an active area of research<sup>1</sup>



<sup>1</sup>G. Ke et al., 'LightGBM: A Highly Efficient Gradient Boosting Decision Tree', *Advances in Neural Information Processing Systems*, ed. by I. Guyon et al., vol. 30, Curran Associates, Inc., 2017

# How to train such models?

Theory

## Ideal scenario

- ★ Split data into training, validation, and test data sets
- ★ Fit model on *training* data, using the *validation* data for hyperparameter tuning, and finally evaluate on *test* data
- ★ Look at *test* data only once!

## Data-sparse scenario

- ★ Follow *cross-validation* procedure, which ensures that every split of the data is used exactly once for testing
- ★ Repeat this procedure multiple times to obtain estimates of generalisation performance
- ★ Can also be combined with scenario above

# How to train such models?

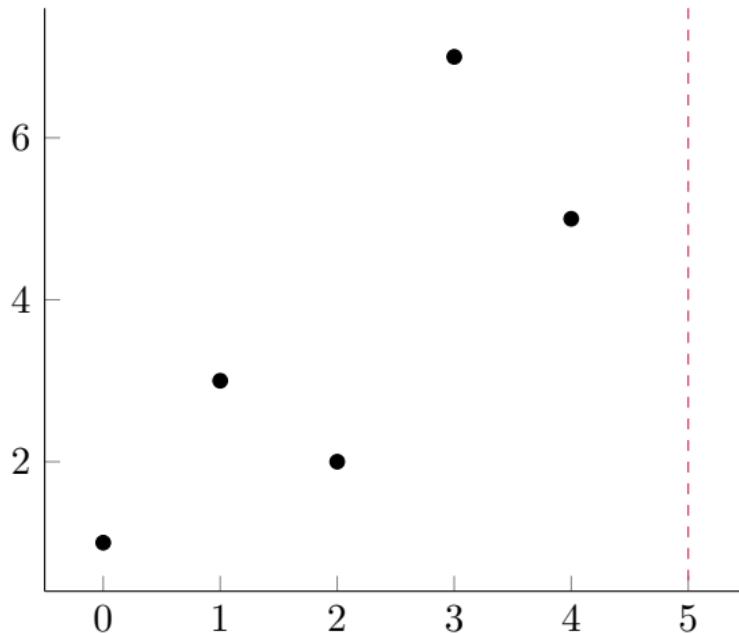
Practice

```
1      cv = StratifiedKFold(n_splits=n_folds, shuffle=True, random_state=42)
2      for train_index, test_index in cv.split(X, y):
3          X_train, y_train = X[train_index], y[train_index]
4          X_test, y_test = X[test_index], y[test_index]
5
6          clf = LogisticRegression(
7              solver="saga",
8              class_weight="balanced",
9              max_iter=5000,
10             random_state=42
11         )
12
13         grid_search = GridSearchCV(
14             clf,
15             param_grid=param_grid,
16             cv=n_folds,
17             scoring="accuracy",
18             n_jobs=-1,
19         )
20
21         grid_search.fit(X_train, y_train)
22         y_pred = grid_search.predict(X_test)
23         accuracy = accuracy_score(y_test, y_pred)
24
```

# Gaussian Processes

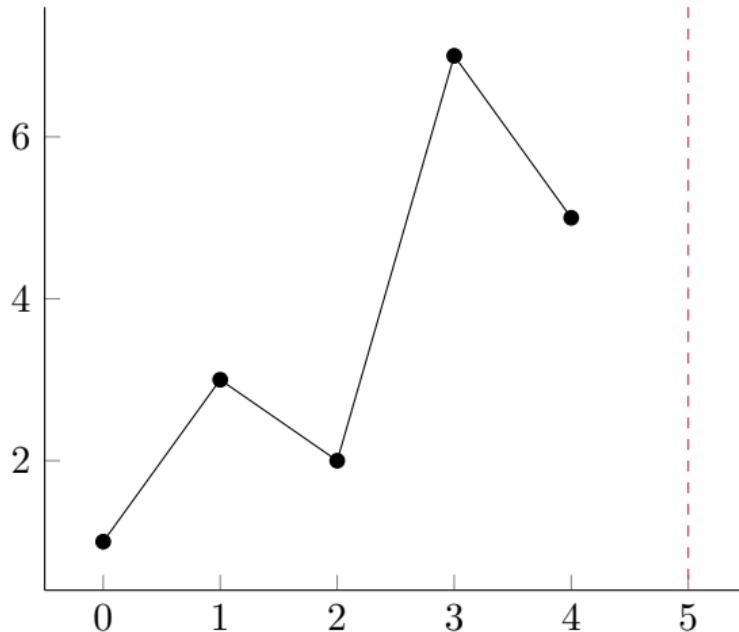
# Motivation

How to interpolate these points?



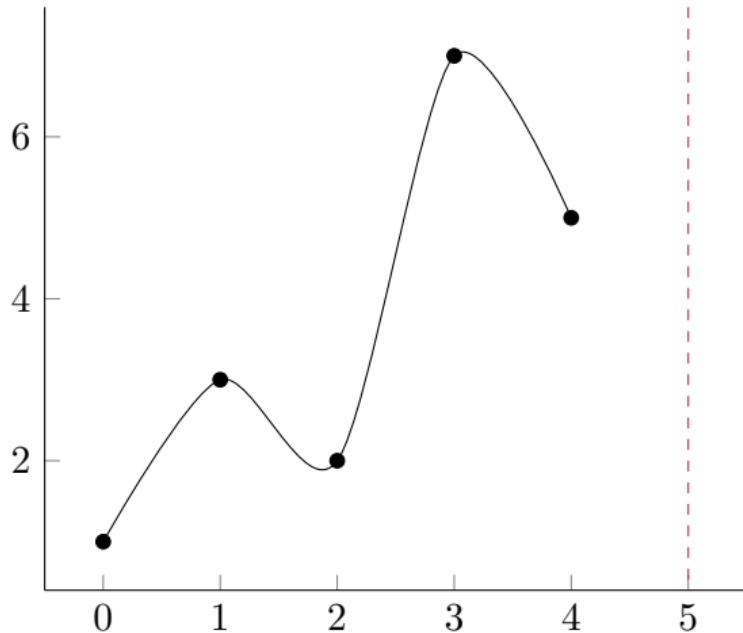
# Motivation

How to interpolate these points?



# Motivation

How to interpolate these points?



# Questions

- ☆ Given known data, what is a *likely* value for an unknown  $x$  coordinate?
- ☆ Can we determine the *confidence* of our prediction?
- ☆ What is the *mean* of all possible predictions?

# Main building block

Multivariate Gaussian distributions

$$f(x_1, \dots, x_k \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{\exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)}{\sqrt{(2\pi)^k \det \boldsymbol{\Sigma}}}$$

- ☆  $\mathbf{x}$  is a  $k$ -dimensional (column) vector
- ☆  $\boldsymbol{\mu}$  is a  $k$ -dimensional (column) vector representing the *mean* of every variable
- ☆  $\boldsymbol{\Sigma}$  is a  $k \times k$  *covariance* matrix
- ☆ One often writes  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

## Conjugate distribution

If  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$  is *jointly Gaussian*, then the *marginals* and the *posterior* are also Gaussian!

# Gaussian processes

C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning*, Adaptive Computation and Machine Learning series, MIT Press, 2005

## Definition

A *Gaussian process* generalises a multivariate normal distribution to *infinitely many* variables. The assumption is that every finite number of variables follows some Gaussian distribution.

## Properties

Gaussian processes can work with different *kernels*, which serve as covariance functions here.

- ☆  $k(x, x') = \sigma_f^2 \exp\left(-\frac{(x-x')^2}{2l^2}\right)$  (*squared exponential kernel*)
- ☆  $k(x, x') = \sigma_b^2 + \sigma_f^2(x - c)(x' - c)$  (*linear kernel*)

# How to perform the regression

Given known values for the function  $f$ , we want to predict new values  $f_*$ . We thus have

$$\begin{pmatrix} f \\ f_* \end{pmatrix} = \mathcal{N} \left( \begin{pmatrix} \mu \\ \mu_* \end{pmatrix}, \begin{pmatrix} \mathbf{K} & \mathbf{K}_* \\ \mathbf{K}_*^T & \mathbf{K}_{**} \end{pmatrix} \right),$$

where we write  $\mathbf{K}_*$  to denote the matrix containing all  $k(x, x_*)$  (and similarly for the other matrices). According to the rules of Gaussian distributions, the *posterior* is given by

$$p(f_* \mid f) = \mathcal{N}(f_* \mid \mu_*, \Sigma_*),$$

where:

$$\mu_* = \mu(\mathbf{x}_*) + \mathbf{K}_*^T \mathbf{K}^{-1} (f - \mu(\mathbf{x}))$$

$$\Sigma_* = \mathbf{K}_{**} - \mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{K}_*$$

# Parameter estimation

We can estimate the parameters that make our observed data most likely. Suppose we have  $k(x, x') = \sigma_f^2 \exp\left(-\frac{(x-x')^2}{2l^2}\right)$ , we can find good values for  $\sigma_f$  and  $l$  by maximising *marginal likelihood*:

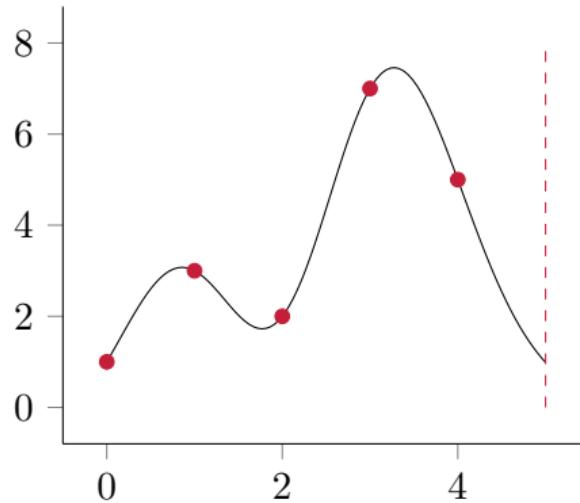
$$\log p(\mathbf{y} | \mathbf{X}) = -\frac{1}{2}\mathbf{y}\mathbf{K}_y^{-1}\mathbf{y} - \frac{1}{2}\log \det \mathbf{K}_y - \frac{n}{2}\log(2\pi)$$

This is possible using *gradient descent* or *Nelder–Mead*, for example.

# How to use this?

Sampling from the *posterior*

Recall that  $p(f_* \mid f) = \mathcal{N}(f_* \mid \mu_*, \Sigma_*)$ , where  $\mu_* = \mu(\mathbf{x}_*) + \mathbf{K}_*^T \mathbf{K}^{-1} (f - \mu(\mathbf{x}))$  and  $\Sigma_* = \mathbf{K}_{**} - \mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{K}_*$ . In our simple example, let's assume that  $\mu = 0$ , so everything reduces to evaluating some matrices.

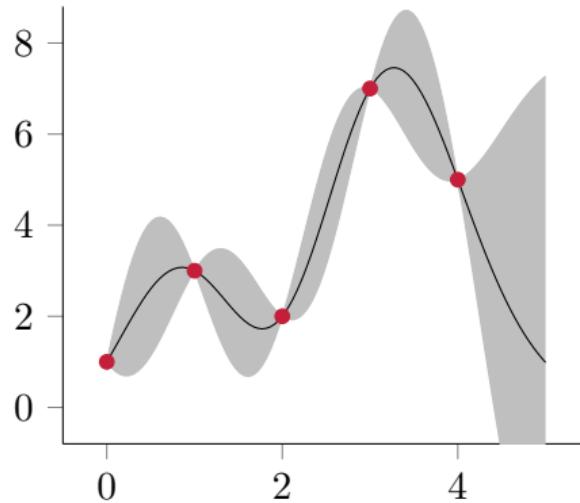


Notice that this is still a *random* sample!

# How to use this?

Sampling from the *posterior*

Recall that  $p(f_* \mid f) = \mathcal{N}(f_* \mid \mu_*, \Sigma_*)$ , where  $\mu_* = \mu(\mathbf{x}_*) + \mathbf{K}_*^T \mathbf{K}^{-1} (f - \mu(\mathbf{x}))$  and  $\Sigma_* = \mathbf{K}_{**} - \mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{K}_*$ . In our simple example, let's assume that  $\mu = 0$ , so everything reduces to evaluating some matrices.



Notice that this is still a *random* sample!

Why kernels are *still* relevant

# Recap

## Kernel

Given a set  $\mathcal{X}$ , a function  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a *kernel* if there is a Hilbert space  $\mathcal{H}$  (an inner product space that is also a complete metric space) and a map  $\Phi: \mathcal{X} \rightarrow \mathcal{H}$ , such that  $k(x, y) = \langle \Phi(x), \Phi(y) \rangle_{\mathcal{H}}$  for all  $x, y \in \mathcal{X}$ .

## Properties

Model similarity between structured and unstructured objects. Many existing machine learning algorithms can be rephrased in terms of kernels using the *kernel trick*.<sup>2</sup>

<sup>2</sup>B. Schölkopf and A. J. Smola, *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*, Adaptive Computation and Machine Learning series, MIT Press, 2018

# Maximum mean discrepancy (MMD)

Comparing distributions of objects using kernels

Given  $n$  samples  $X = \{x_1, \dots, x_n\} \subseteq \mathcal{X}$  and  $m$  samples  $Y = \{y_1, \dots, y_m\} \subseteq \mathcal{X}$ , the biased empirical estimate of the MMD between  $X$  and  $Y$  is

$$\text{MMD}^2(X, Y) := \frac{1}{n^2} \sum_{i,j=1}^n k(x_i, x_j) + \frac{1}{m^2} \sum_{i,j=1}^m k(y_i, y_j) - \frac{2}{nm} \sum_{i=1}^n \sum_{j=1}^m k(x_i, y_j).$$

This formulation is due to Gretton et al.<sup>3</sup> under some mild conditions, MMD is a metric on the space of probability distributions.

<sup>3</sup>A. Gretton, K. Borgwardt, M. Rasch, B. Schölkopf and A. J. Smola, 'A Kernel Method for the Two-Sample-Problem', *Advances in Neural Information Processing Systems*, ed. by B. Schölkopf, J. C. Platt and T. Hoffman, vol. 19, MIT Press, 2007, pp. 513–520

# Deep Learning

# Some terminology

- ☆ *Batch* or *Mini-batch*: a subsample of the data set
- ☆ *Epoch*: how many times the entire data set was processed
- ☆ *Loss*: a function that determines how well the network fits the data

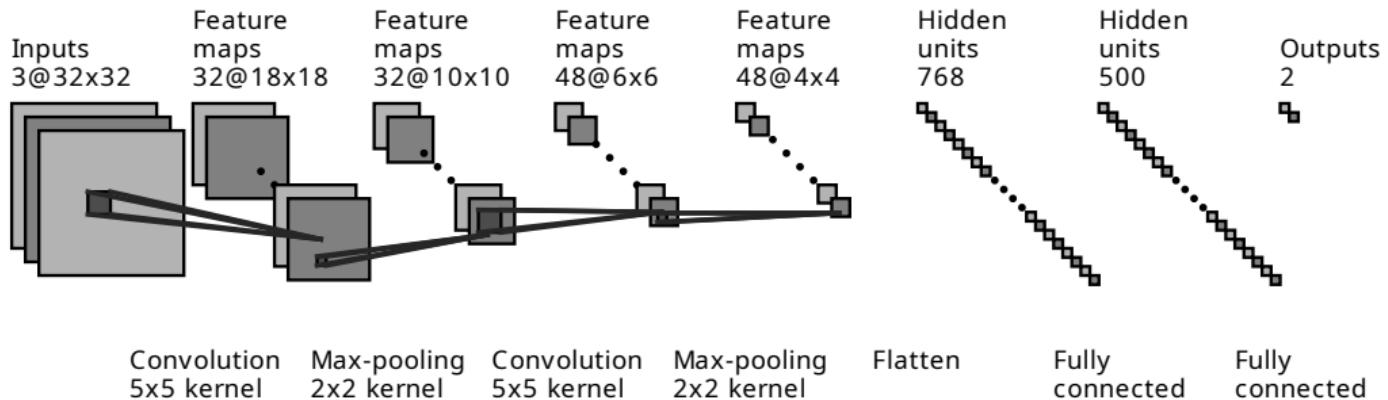
## Backpropagation

Adjust weights—connection strengths—in the neural network to minimise the error: compute gradient of the loss function for this batch and update weights.

# Some common architectures

- ☆ Fully-connected neural networks (FCNs)
- ☆ Convolutional neural networks (CNNs)
- ☆ Autoencoders (AEs)
- ☆ Recurrent neural networks (RNNs)

# Convolutional neural networks



# The ecosystem of deep learning

- ☆ Keras
- ☆ PyTorch
- ☆ TensorFlow

# Example

Getting some data

```
1 import torch
2 from torch import nn
3 from torch.utils.data import DataLoader
4 from torchvision import datasets
5 from torchvision.transforms import ToTensor
6
7 training_data = datasets.FashionMNIST(
8     root="data", train=True, download=True,
9     transform=ToTensor()
10)
11
12 test_data = datasets.FashionMNIST(
13     root="data", train=False, download=True,
14     transform=ToTensor()
15)
16
17 train_dataloader = DataLoader(training_data, batch_size=64)
18 test_dataloader = DataLoader(test_data, batch_size=64)
19
```

(pytorch.org)

# Example

Defining a neural network

```
1  class NeuralNetwork(nn.Module):
2      def __init__(self):
3          super(NeuralNetwork, self).__init__()
4          self.flatten = nn.Flatten()
5          self.linear_relu_stack = nn.Sequential(
6              nn.Linear(28*28, 512),
7              # Non-linear operation!
8              nn.ReLU(),
9              nn.Linear(512, 512),
10             nn.ReLU(),
11             nn.Linear(512, 10),
12         )
13
14     def forward(self, x):
15         x = self.flatten(x)
16         logits = self.linear_relu_stack(x)
17         return logits
18
19 model = NeuralNetwork()
20
```

(pytorch.org)

# Example

## Train loop

```
1 def train_loop(dataloader, model, loss_fn, optimizer):
2     size = len(dataloader.dataset)
3     for batch, (X, y) in enumerate(dataloader):
4         # Compute prediction and loss
5         pred = model(X)
6         loss = loss_fn(pred, y)
7
8         # Backpropagation
9         optimizer.zero_grad()
10        loss.backward()
11        optimizer.step()
12
```

(pytorch.org)

# Example

## Test loop

```
1 def test_loop(dataloader, model, loss_fn):
2     size = len(dataloader.dataset)
3     num_batches = len(dataloader)
4     test_loss, correct = 0, 0
5
6     with torch.no_grad():
7         for X, y in dataloader:
8             pred = model(X)
9             test_loss += loss_fn(pred, y).item()
10            correct += (pred.argmax(1) == y).type(torch.float).sum().item()
11
12    test_loss /= num_batches
13    correct /= size
14
```

(pytorch.org)

# Example

Running the training and testing process

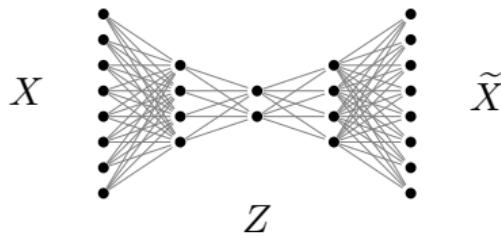
```
1 learning_rate = 1e-3
2 batch_size = 64
3 epochs = 10
4
5 loss_fn = nn.CrossEntropyLoss()
6 optimizer = torch.optim.SGD(model.parameters(), lr=learning_rate)
7
8 for t in range(epochs):
9     train_loop(train_dataloader, model, loss_fn, optimizer)
10    test_loop(test_dataloader, model, loss_fn)
11
```

(pytorch.org)

# A lot of moving parts

- ☆ Architecture
- ☆ Loss function
- ☆ Optimiser
- ☆ Additional elements: batch norm, attention, ...

# Autoencoders



## Terminology

- ☆  $X \in \mathbb{R}^D$ : input data
- ☆  $Z \in \mathbb{R}^d$ : latent representation
- ☆  $\tilde{X} \in \mathbb{R}^D$ : reconstructed data

## Properties

- ☆ Typically,  $D \gg d$ .
- ☆ Use loss function  $\mathcal{L}(X, \tilde{X})$  to measure quality of reconstruction.

# Why autoencoders?

- ☆ Encoder–decoder architecture (we learn the *identity* function).
- ☆ ‘Middle’ layer serves as *bottleneck* or *latent representation*.
- ☆ Latent representations can be used for visualisation, interpolation, clustering, and much more.

# A simple autoencoder

- ☆ Encoder: linear transformation  $\mathbb{R}^D \rightarrow \mathbb{R}^2$
- ☆ Decoder: linear transformation  $\mathbb{R}^2 \rightarrow \mathbb{R}^D$
- ☆ Loss function: mean squared error,  $\mathcal{L}(X, \tilde{X}) := \|X - \tilde{X}\|_2^2$

# A simple autoencoder

Some reconstructions



o epochs

# A simple autoencoder

Some reconstructions



10 epochs

# A simple autoencoder

Some reconstructions



20 epochs

# A simple autoencoder

Some reconstructions



30 epochs

# A simple autoencoder

Some reconstructions



40 epochs

# A simple autoencoder

Some reconstructions



50 epochs

# PyTorch Lightning

## Overall training loop

```
1 from .data import MNISTDataModule
2 from .models import LinearAutoencoder
3 from pytorch_lightning import Trainer
4
5 if __name__ == '__main__':
6     dm = MNISTDataModule()
7     dm.prepare_data()
8     dm.setup()
9
10    model = LinearAutoencoder(input_dim=dm.input_dim, bottleneck_dim=2)
11
12    # There's a lot more parameters for a 'trainer' class, but this is
13    # sufficient for our brief example.
14    trainer = Trainer(max_epochs=50)
15    trainer.fit(model, dm)
16
```

# PyTorch Lightning

## Defining a model

```
1 import torch.nn as nn
2 import torch.nn.functional as F
3 import pytorch_lightning as pl
4
5 class LinearAutoencoder(pl.LightningModule):
6     def __init__(self, input_dim, bottleneck_dim=2, lr=0.01):
7         super().__init__()
8
9         self.encoder = nn.Linear(input_dim, bottleneck_dim)
10        self.decoder = nn.Linear(bottleneck_dim, input_dim)
11
12        self.lr = 0.01
13        self.loss_fn = F.mse_loss
14
```

# PyTorch Lightning

Defining a model, continued

```
1  def forward(self, x):
2      z = self.encoder(x)
3      x_hat = self.decoder(z)
4      loss = self.loss_fn(x_hat, x)
5      return x_hat, z, loss
6
7  def training_step(self, batch, batch_idx):
8      x, y = batch
9
10     x = x.view(x.size(0), -1)
11     _, loss = self(x)
12
13     self.log('train_loss', loss)
14     return loss
15
16  def configure_optimizers(self):
17      optimizer = torch.optim.Adam(self.parameters(), lr=self.lr)
18      return optimizer
19
```

# PyTorch Lightning

More complex operations

## Some highlights

- ☆ Early stopping
- ☆ Learning rate adjustments
- ☆ Model checkpoints
- ☆ Distributed training

# Variational autoencoders

We can make the autoencoder learn the *parameters* of our input distribution to sample *new* data points from it.

# Variational autoencoders

The basic skeleton

```
1  class BetaVAE(pl.LightningModule):
2      def __init__(self, input_dim, bottleneck_dim=32, beta=1, lr=1e-3):
3          super(BetaVAE, self).__init__()
4
5          self.bottleneck_dim = bottleneck_dim
6
7          # The specific choice of encoder/decoder is irrelevant here.
8          self.encoder = ...
9          self.decoder = ...
10
11         # Final layer for encoding the parameters of the distribution
12         self.fc_mu = nn.Linear(16, bottleneck_dim)
13         self.fc_logvar = nn.Linear(16, bottleneck_dim)
14
15         # For sampling values that parametrise our distribution
16         self.fc_z = nn.Linear(bottleneck_dim, 16)
17
```

# Variational autoencoders

Encoding and decoding

```
1  def encode(self, x):
2      x = self.encoder(x)
3      x = x.view(x.shape[0], -1)
4      return self.fc_mu(x), self.fc_logvar(x)
5
6  def sample(self, mu, logvar):
7      std = torch.exp(0.5 * logvar)    # e^(1/2 * log(std^2))
8      eps = torch.randn_like(std)    # random ~ N(0, 1)
9      return eps.mul(std).add_(mu)
10
11 def decode(self, z):
12     z = self.fc_z(z)
13     z = z.view(z.shape[0], -1)
14     return self.decoder(z)
15
16 def forward(self, x):
17     mu, logvar = self.encode(x)
18     z = self.sample(mu, logvar)
19     x_hat = self.decode(z)
20     return x_hat, mu, self.loss_fn(x_hat, x, mu, logvar)
21
```

# Variational autoencoders

Loss function

```
1     def loss_fn(self, x_hat, x, mu, logvar):
2         reconstruction_loss = F.mse_loss(
3             x_hat,
4             x,
5             reduction='sum' # We reduce it ourselves later on!
6         )
7
8         # This follows the loss provided by Kingma and Welling in
9         # 'Auto-Encoding Variational Bayes', Appendix B. The loss
10        # assumes that all distributions are Gaussians.
11        kl = -0.5 * torch.sum(1 + logvar - mu.pow(2) - logvar.exp())
12
13        return (reconstruction_loss + self.beta * kl) / x.shape[0]
```

The loss term boils down to a Kullback–Leibler divergence between the prior distribution and the posterior distribution:

$$\text{KL}(q_\phi(\mathbf{z}) \parallel p_\theta(\mathbf{z}))$$

# Variational autoencoders

Some reconstructions



o epochs

# Variational autoencoders

Some reconstructions



10 epochs

# Variational autoencoders

Some reconstructions



20 epochs

# Variational autoencoders

Some reconstructions



30 epochs

# Variational autoencoders

Some reconstructions



40 epochs

# Variational autoencoders

Some reconstructions



50 epochs

# Sampling from the learned distribution

Pick  $z \sim \mathcal{N}(0, 1)$  and *decode* it. This yields a valid image of the distribution!



# Brief summary

- ☆ Taming the dragon: PyTorch Lightning makes working with neural networks easy
- ☆ Autoencoders are available in many flavours (variational, topological, ...)
- ☆ *Acta non verba*: try them out yourselves!

<https://github.com/Pseudomanifold/bestiary-autoencoders>