

---

# Lecture Recording

---

- ❖ **Note: These lectures will be recorded and posted onto the IMPRS website**
- ❖ Dear participants,
- ❖ We will record all lectures on “*Making sense of data: introduction to statistics for gravitational wave astronomy*”, including possible Q&A after the presentation, and we will make the recordings publicly available on the IMPRS lecture website at:
  - <https://imprs-gw-lectures.aei.mpg.de/2023-making-sense-of-data/>
- ❖ By participating in this Zoom meeting, you are giving your explicit consent to the recording of the lecture and the publication of the recording on the course website.

# Making sense of data: introduction to statistics for gravitational wave astronomy

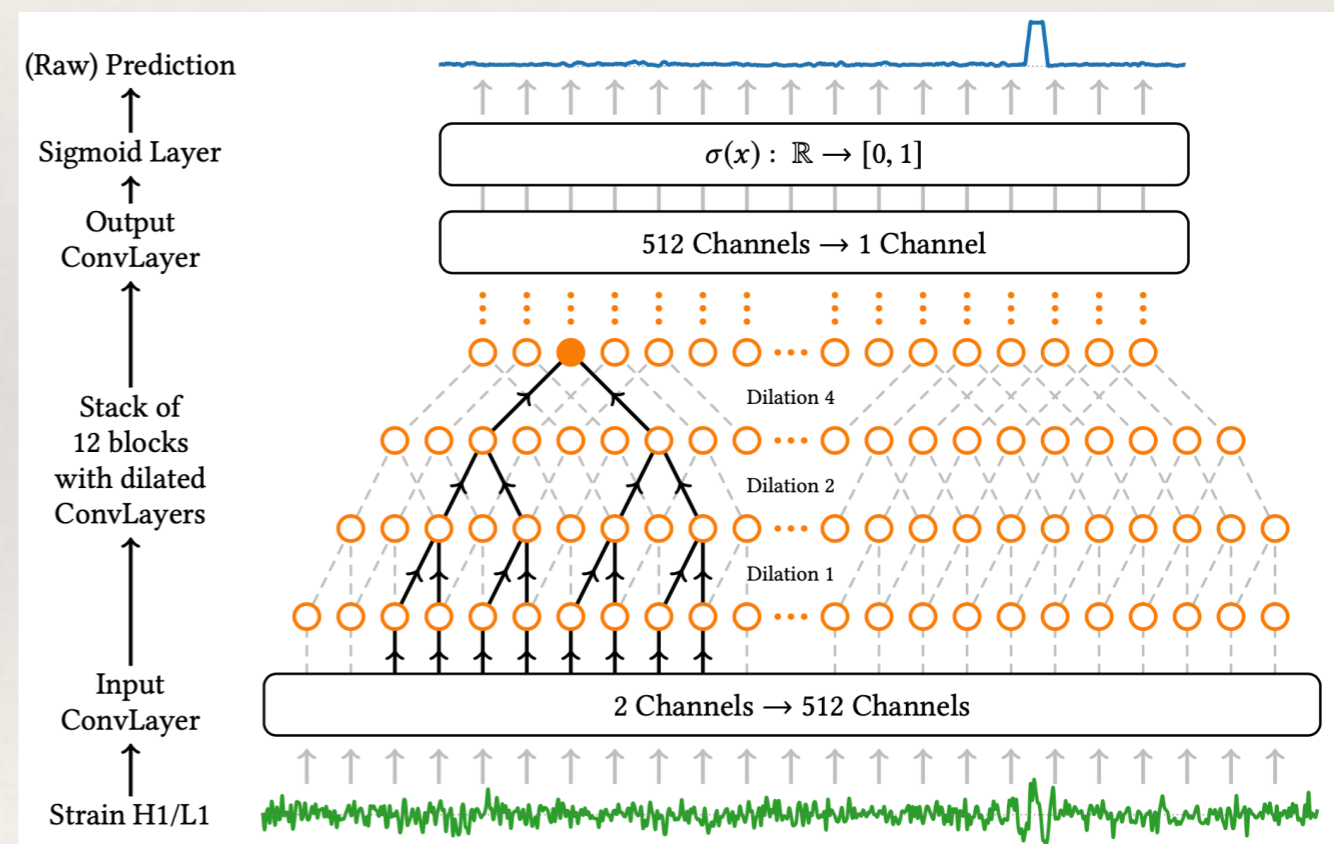
## Part III: Machine Learning

### Lecture 1: Introduction

AEI IMPRS Lecture Course

Jonathan Gair [jgair@aei.mpg.de](mailto:jgair@aei.mpg.de)

Thanks to **Stephen Green** for producing much of the material for this course in 2021!



---

# Course outline

---

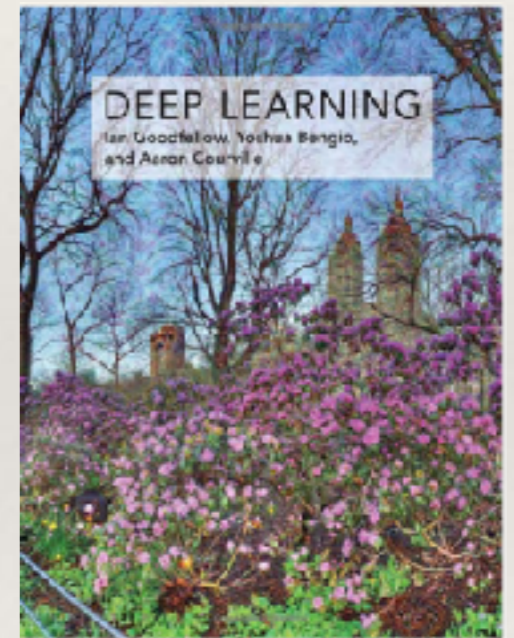
- ❖ Lecture 1: introduction to machine learning
- ❖ Lecture 2: neural networks and deep learning
- ❖ Lecture 3: machine learning for gravitational wave astronomy
- ❖ Practical: GW search and parameter estimation using machine learning

---

# References

---

- ❖ **Textbook: “Deep Learning” by Goodfellow, Bengio, and Courville**
  - Free online at <https://www.deeplearningbook.org>
  - Course covers parts of Chapters 5, 6, 9, 20
  
- ❖ **pyTorch**
  - machine learning framework for practical part
  - many tutorials at <https://pytorch.org>



---

# Introduction to machine learning

---

- ❖ Computers are designed to complete repetitive tasks. A task typically involves taking an *input* and mapping it to an *output*.
- ❖ A computer programme is a *set of instructions* that *teach* the computer how to perform a task.
- ❖ *Machine learning* is the development of approaches that allow computers to *learn* how to perform a task, typically by seeing a large set of examples.
- ❖ Machine learning algorithms typically consist of *function approximators* that have a *large number of free parameters*. These are designed in a way that allow the choice of parameters to be *automatically optimised* to *minimise* a specified *objective function* (the loss function).

# Introduction to machine learning

- ❖ **Example**

- ❖ *Classification*: learn a function that maps input data into a category

$$f : \mathbb{R}^n \rightarrow \{1, \dots, k\}$$

- ❖ e.g., recognise handwriting digits



8, 8, 1, 5, 1  
4, 4, 7, 4, 9

---

# Machine Learning Tasks

---

❖ Examples:

- **Regression:** Learn a function predicting real-valued quantities

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

E.g., What are the physical parameters characterizing a binary merger?

- **Sampling:** Generate new samples similar to training examples.
- **Denoising:** Given noisy data  $\tilde{\mathbf{x}} \in \mathbb{R}^n$ , predict clean data  $\mathbf{x} \in \mathbb{R}^n$ :  $p(\mathbf{x} | \tilde{\mathbf{x}})$
- **Density estimation:** Given training examples  $\mathbf{x} \in \mathbb{R}^n$  learn a probability density function  $p(\mathbf{x})$ .
- **Game playing:** Given a game configuration, what is the best move to make?

---

# Performance Measures

---

- ❖ For each task, it is necessary to specify some quantitative measure of performance:
  - for **classification**, the accuracy (the fraction of examples that produce the correct output)
  - for **density estimation**, the log probability assigned to examples
  - for **regression**, the mean squared error
- ❖ We are usually interested in how the machine learning algorithm performs on data that have not been seen before: Evaluate performance on a **test set** that is different from the **training set**.



---

# Types of Learning Algorithms

---

- ❖ Typically we have a dataset  $\{\mathbf{x}^{(i)}\}$  consisting of many data points  $\mathbf{x}^{(i)} \in \mathbb{R}^n$ . The data points may or may not have associated labels  $\mathbf{y}^{(i)} \in \mathbb{R}^m$ .
  - ❖ **Unsupervised:** learn  $p(\mathbf{x})$ 
    - Examples: density estimation, sampling
  - ❖ **Supervised:** learn  $p(\mathbf{y} | \mathbf{x})$ 
    - Examples: regression, classification
- ❖ **Reinforcement** learning allows the algorithm to interact with the environment and produce new samples (e.g., game playing).

somewhat hazy distinction,  
e.g., learning  $p(\mathbf{y}, \mathbf{x})$

---

# Maximum likelihood estimation

---

- ❖ Consider a set of  $N$  independent examples  $\mathbf{x}^{(i)} \sim p_{\text{data}}(\mathbf{x})$  drawn from the data-generating distribution.
- ❖ **Unsupervised learning:** Let  $p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$  be a parametric family of model probability distributions. Choose  $\boldsymbol{\theta}$  such that this becomes a good approximation to  $p_{\text{data}}(\mathbf{x})$ .
- ❖ **Maximum likelihood estimator** is 
$$\begin{aligned}\boldsymbol{\theta}_{\text{ML}} &= \arg \max_{\boldsymbol{\theta}} p_{\text{model}}(\mathbf{X}; \boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^N p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \log p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})\end{aligned}$$
- ❖ Equivalent to minimizing **KL divergence** or **cross-entropy** between  $p_{\text{data}}$  and  $p_{\text{model}}$ .

---

# Conditional Estimation

---

- ❖ **Supervised learning:** Estimate a conditional probability  $p_{\text{model}}(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta})$
- ❖ Generalize the maximum likelihood estimator:

$$\begin{aligned}\boldsymbol{\theta}_{\text{ML}} &= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log p_{\text{model}}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{p_{\text{data}}(\mathbf{x}, \mathbf{y})} \log p_{\text{model}}(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta})\end{aligned}$$

- ❖ This is one of the most common situations.

# Example: Linear regression

- ❖ Suppose we have labelled data  $(\mathbf{x}^{(i)}, y^{(i)})$ .
- ❖ Let  $p(y|\mathbf{x}) = \mathcal{N}(\mu(\mathbf{x}), \sigma^2)(y)$  where  $\mu(\mathbf{x}) = \boldsymbol{\theta} \cdot \mathbf{x}$ ;  $\sigma$  fixed.

- ❖ Using the PDF  $p(y|\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \mu(\mathbf{x}))^2}{2\sigma^2}\right)$

we obtain the loss function

$$J(\boldsymbol{\theta}) = -\sum_{i=1}^N \log p(y^{(i)}|\mathbf{x}^{(i)}; \boldsymbol{\theta}) \quad \propto \text{mean squared error}$$

$$= \frac{N}{2} \log 2\pi\sigma^2 + \sum_{i=1}^N \frac{(y^{(i)} - \mu(\mathbf{x}^{(i)}))^2}{2\sigma^2}$$

- ❖ Can solve exactly  $\nabla_{\boldsymbol{\theta}} J = 0 \implies \boldsymbol{\theta}_{\text{ML}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$

---

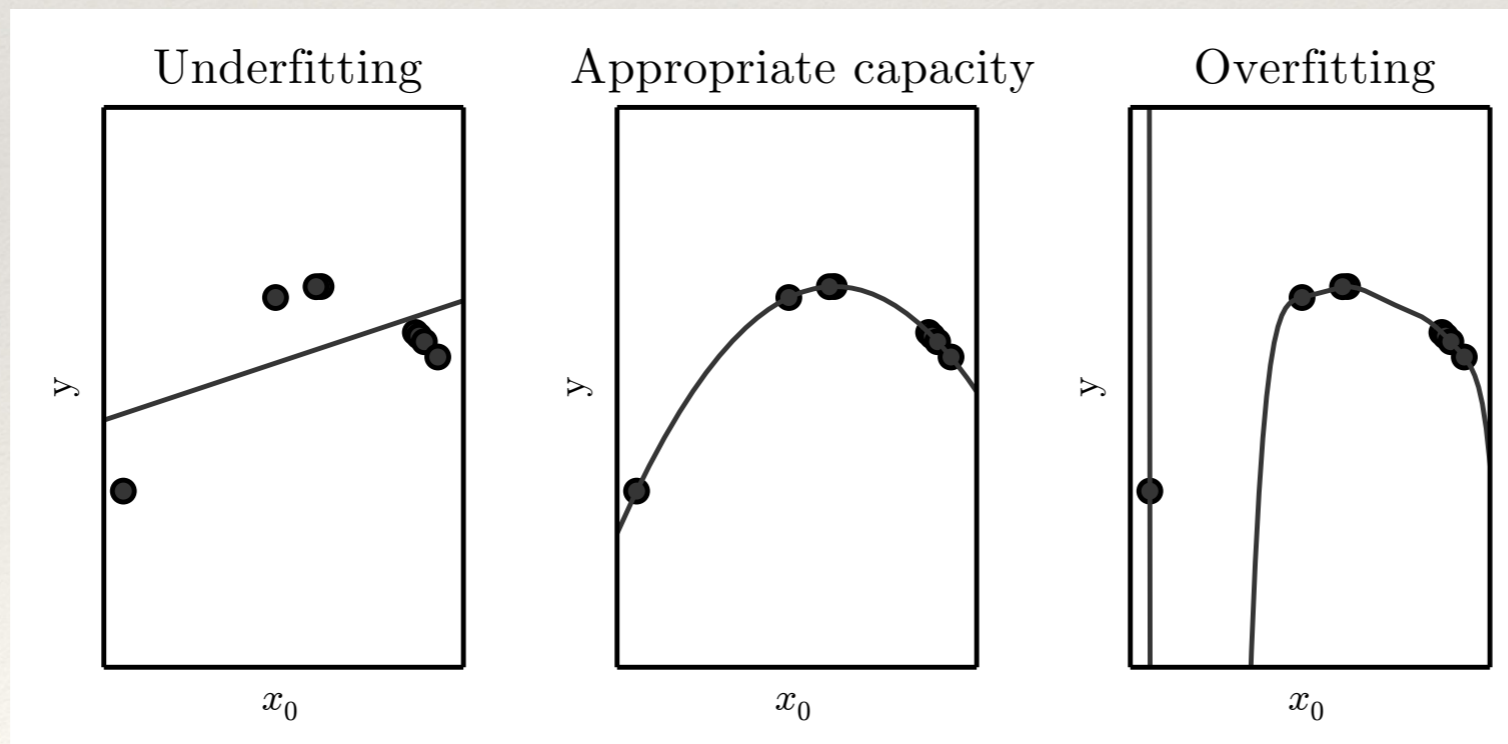
# More general regression

---

- ❖ More generally  $\mu(\mathbf{x})$  does not have to be linear. We can increase the **representational capacity** of the model by using more complicated functions.
  - E.g., polynomial  $\mu(x) = b + \sum_{i=1}^k w_i x^i$  (can still solve in closed form)
    - ← “hyperparameter”
  - E.g. nonparametric regression
    - nearest neighbor: For any  $\mathbf{x}$ , find the nearest  $\mathbf{x}^{(i)}$  in the training set and return  $y^{(i)}$ .
  - E.g., neural network (next lecture)
- ❖ Not all models can be optimized in closed form. The optimization algorithm may be imperfect, so the **effective capacity** is lower than the representational capacity.

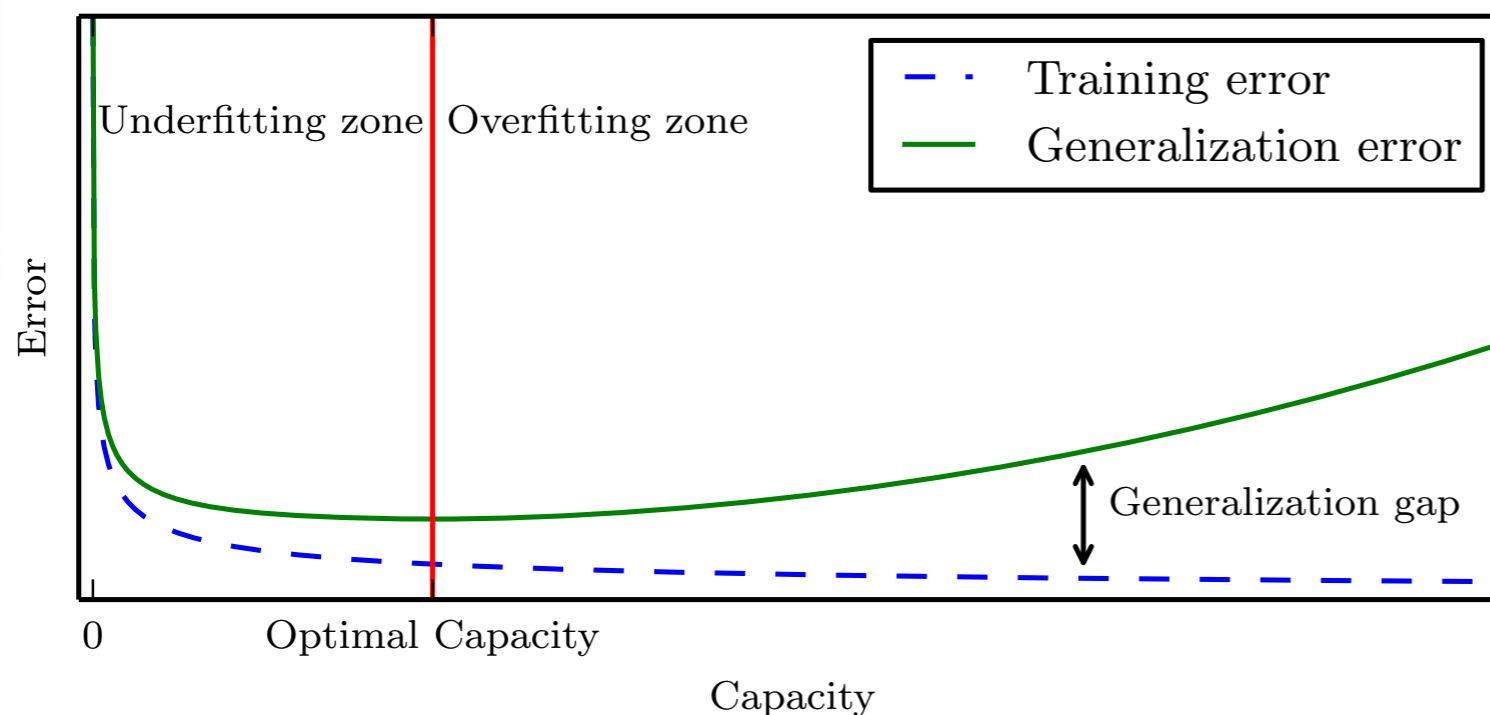
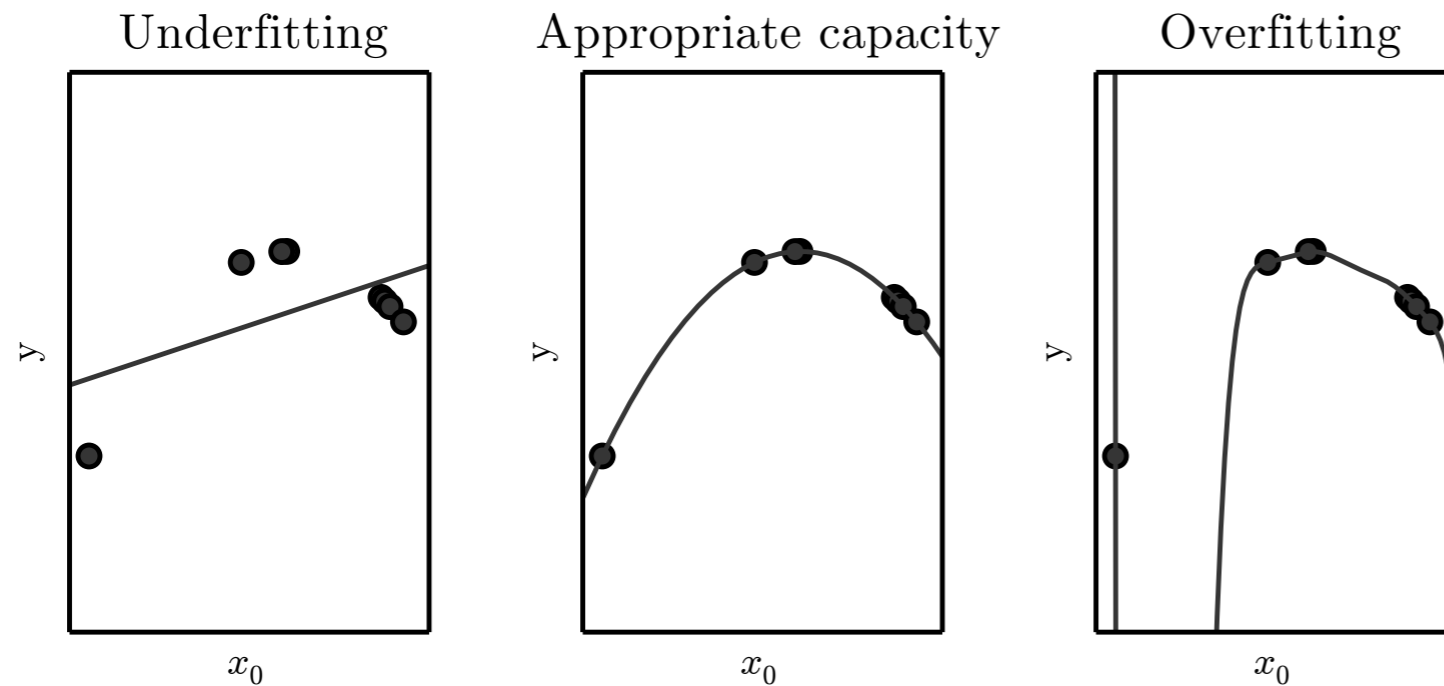
# Overfitting and underfitting

- ❖ Higher capacity models run the risk of **overfitting**. The algorithm must perform well not just on data used for training, but also on new, previously unseen inputs (test data). This is called **generalization**.
- ❖ Training and test examples should be **independent and identically distributed (i.i.d.)**, i.e., drawn from the same data-generating distribution  $p_{\text{data}}$



Goodfellow et al (2016)

# Overfitting and underfitting

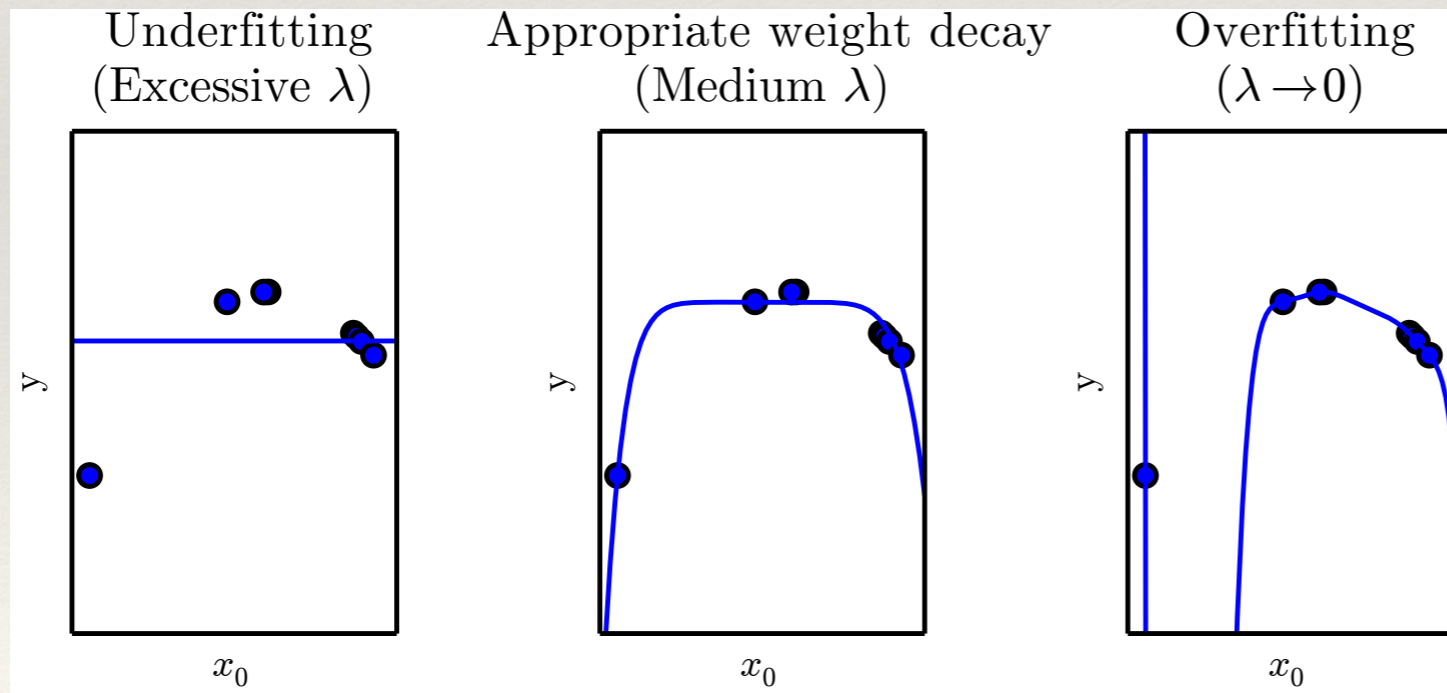


- ❖ Capacity should be chosen to minimize generalization error.
- ❖ Depends also on the size of the training set.

# Regularization

- ❖ One way to improve generalization is to build in preferences for certain values of the parameters  $\theta$ , without changing the representational capacity.
- ❖ Add a **regularizer** to the loss function.

Weight decay:  $J(\theta) = \text{MSE} + \lambda \theta^\top \theta$  preference for small values of  $\theta$



But do we still have a probabilistic interpretation of this loss?

Goodfellow et al (2016)



---

# Bayesian statistics for model parameters

---

- ❖ The maximum likelihood objective picks out a single choice of parameters  $\theta_{\text{ML}}$  corresponding to the maximum of  $p(\mathbf{X} | \theta)$ .
- ❖ We can also treat  $\theta$  in a Bayesian way:
  - Specify a prior  $p(\theta)$
  - Obtain the posterior using Bayes' rule  $p(\theta | \mathbf{X}) = \frac{p(\mathbf{X} | \theta)p(\theta)}{p(\mathbf{X})}$
- ❖ This incorporates the **uncertainty** associated to the choice of  $\theta$ .
- ❖ The **prior** acts as a regularizer.

# Example: Bayesian linear regression

- ❖ As before we take a Gaussian likelihood  $p(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \mathcal{N}(\mathbf{X}\mathbf{w}, I)(\mathbf{y})$
- ❖ Also take a Gaussian prior  $p(\mathbf{w}) = \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Lambda}_0)$
- ❖ Exercise: show that the posterior is also Gaussian, of the form

$$p(\mathbf{w} | \mathbf{X}, \mathbf{y}) \propto \exp \left( -\frac{1}{2} (\mathbf{w} - \boldsymbol{\mu}_m)^\top \boldsymbol{\Lambda}_m^{-1} (\mathbf{w} - \boldsymbol{\mu}_m) \right)$$

$$\boldsymbol{\Lambda}_m = (\mathbf{X}^\top \mathbf{X} + \boldsymbol{\Lambda}_0^{-1})^{-1}$$

$$\boldsymbol{\mu}_m = \boldsymbol{\Lambda}_m (\mathbf{X}^\top \mathbf{y} + \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0)$$

---

# Maximum a posteriori estimation

---

- ❖ To obtain a point estimate that still takes into account prior, we can take the maximum of the posterior distribution over  $\boldsymbol{\theta}$ ,

$$\begin{aligned}\boldsymbol{\theta}_{\text{MAP}} &= \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta}|\mathbf{x}) \\ &= \arg \max_{\boldsymbol{\theta}} (\log p(\mathbf{x}|\boldsymbol{\theta}) + \log p(\boldsymbol{\theta}))\end{aligned}$$

For  $p(\mathbf{w}) = \mathcal{N}(0, \mathbf{I}/\lambda)$  this term  $\longrightarrow \lambda \mathbf{w}^T \mathbf{w}$

- ❖ MAP Bayesian inference with a Gaussian weight prior corresponds to weight decay. More generally, MAP provides a way to interpret regularization terms.

---

# Example: Logistic regression

---

- ❖ If instead of estimating real-valued quantity  $y$ , we are interested in a binary classification problem with  $y \in \{0,1\}$ , we can use **logistic regression**.
- ❖ Use a logistic sigmoid function  $\sigma(u) = \frac{1}{1 + e^{-u}}$  to squeeze the result of linear regression to lie between 0 and 1. Interpret as a probability
$$p(y = 1|\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^\top \mathbf{x})$$
- ❖ Can use maximum likelihood estimation to determine parameters  $\mathbf{w}$ . But there is no analytic solution because of nonlinearity.

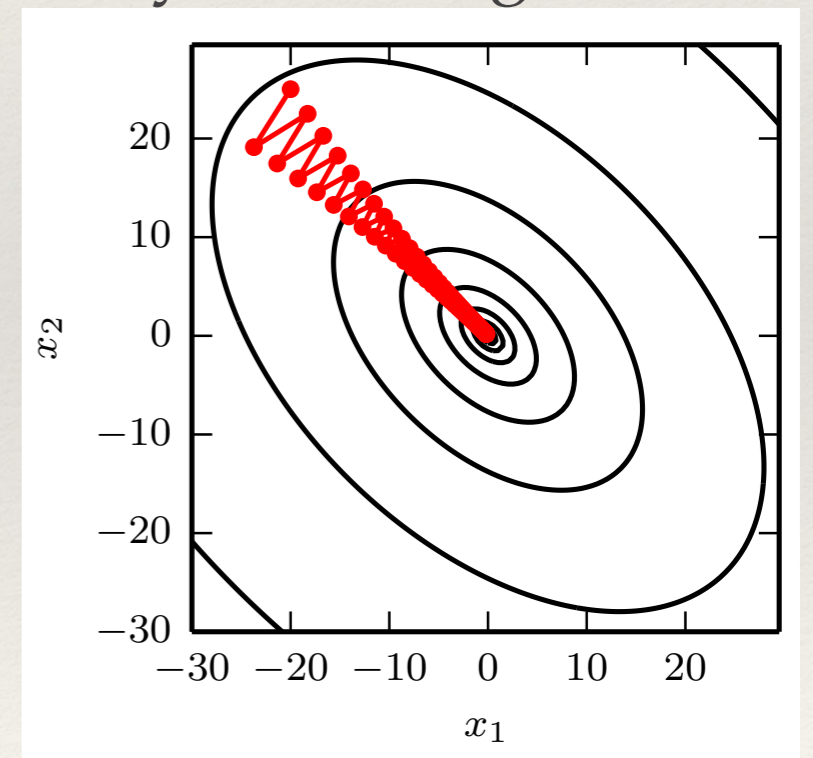
# Stochastic gradient descent

- ❖ In the case where a closed-form minimum is not available, **gradient descent** can be used to **optimize** the loss function, i.e., to tune  $\theta$  to approach the minimum.
- ❖ Starting from a point  $\theta_0$  we can move to a new point by following the gradient

$$\theta_1 = \theta_0 - \epsilon \nabla_{\theta} J |_{\theta_0}$$

“Learning rate”

- ❖ Higher order algorithms can involve the second or higher derivatives (e.g., Hessian).



Goodfellow et al (2016)

---

# Stochastic gradient descent

---

- ❖ For the negative log likelihood loss, the gradient reduces to the sum of per-example gradients,

$$\nabla_{\boldsymbol{\theta}} J = -\frac{1}{N} \sum_{i=1}^N \log p(y^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta})$$

↙ Cost  $\propto N$

- ❖ Can break this up into **minibatches** (subsets of the full training set). Typically this could be several hundred training elements.
- ❖ This has two main advantages: (1) it is faster to compute each update, and (2) it introduces stochasticity, which helps avoid local minima.

---

# Summary

---

- ❖ A machine learning algorithm requires the following:
  1. **dataset** —  $\{\mathbf{x}^{(i)}, y^{(i)}\}$  (supervised) or  $\{\mathbf{x}^{(i)}\}$  (unsupervised)
  2. **model** — E.g., linear regression  $p_{\text{model}}(y | \mathbf{x}) = \mathcal{N}(\boldsymbol{\theta}^\top \mathbf{x}, 1)(y)$
  3. **loss function** — E.g.,  $J(\theta) = -\mathbb{E}_{p_{\text{data}}(\mathbf{x})} \log p_{\text{model}}(\mathbf{x})$
  4. **optimization algorithm** — E.g., stochastic gradient descent

---

# Next lecture: deep learning

---

## ❖ Challenges:

- High dimensionality of data:

The number of possible data configurations is exponential in the number of data dimensions. Hard to cover this with training data.

- Manifold learning:

For many data sets, actual data realizations form a much lower dimensional subset of  $\mathbb{R}^n$ . E.g., random realizations of images will look like noise.