

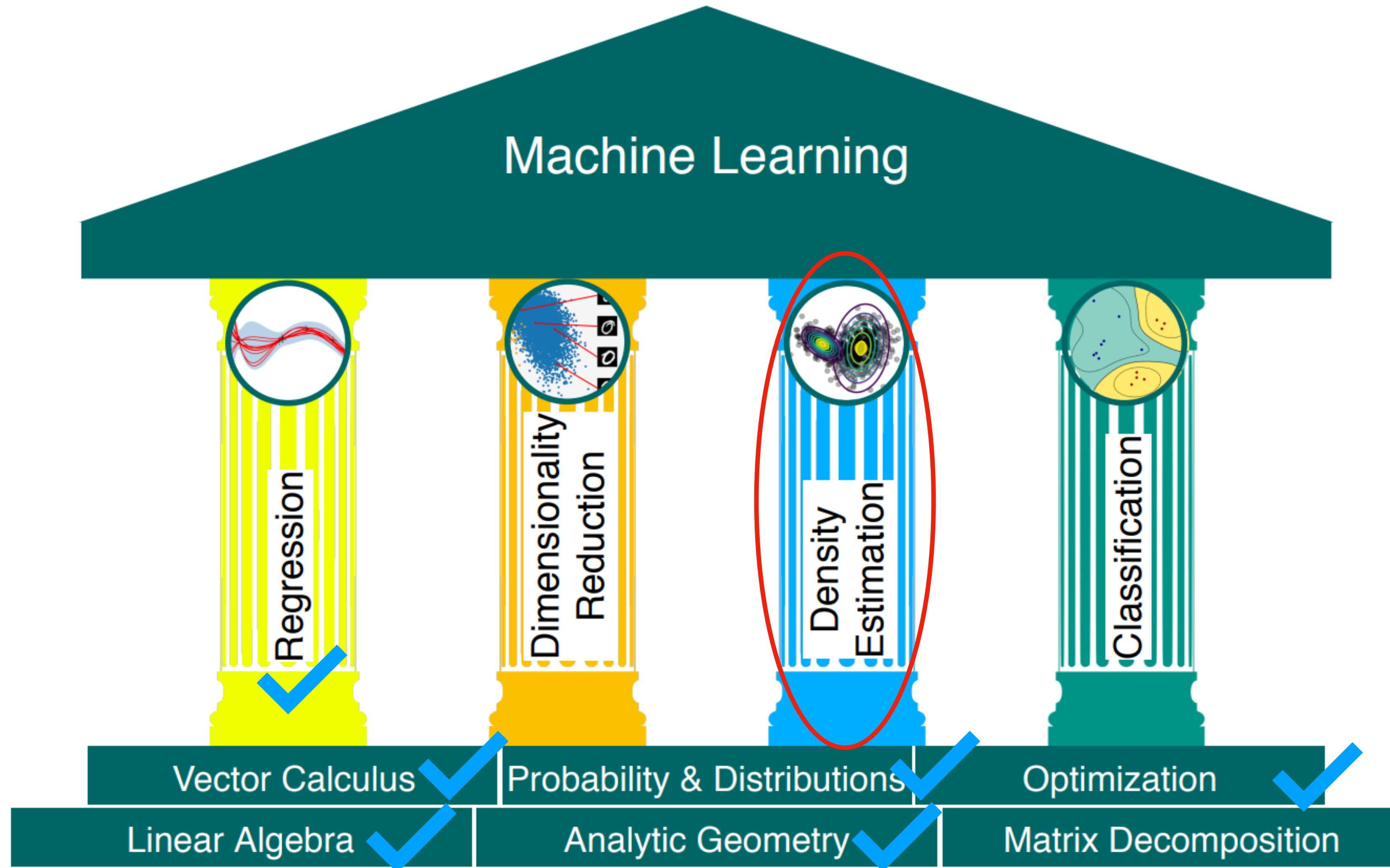
# **Density estimation with Gaussian mixture models**

**Week 8 - Introduction to ML / Thang Bui / ANU / 2023 S2**

# Housekeeping

- Assignment 3 is now available
- [Centrally invigilated, written] exam timetable will be available soon
- Please use Ed!
- Next Monday is a public holiday
  - Lecture: Dunbar Lecture Theatre, Physics Bldg 39A, Thursday 8:30 - 10 am
  - Tutorial: Comp Lab 1.24, Hanna Neumann Bldg 145, Tuesday 5 pm

# Foundations of ML



This week:

Density estimation using Gaussian mixture models

# Recap: Linear regression

Point estimate

$$\theta_{\text{ML}} = (X^T X)^{-1} X^T \mathbf{y} \text{ or } \theta_{\text{MAP}} = (X^T X + N\lambda \mathbf{I})^{-1} X^T \mathbf{y}$$

Exact posterior

$$p(\theta | \mathcal{D}) = \mathcal{N}(\theta; \mu, \Sigma)$$

$$\mu = \Sigma \sigma^{-2} X^T \mathbf{y}$$

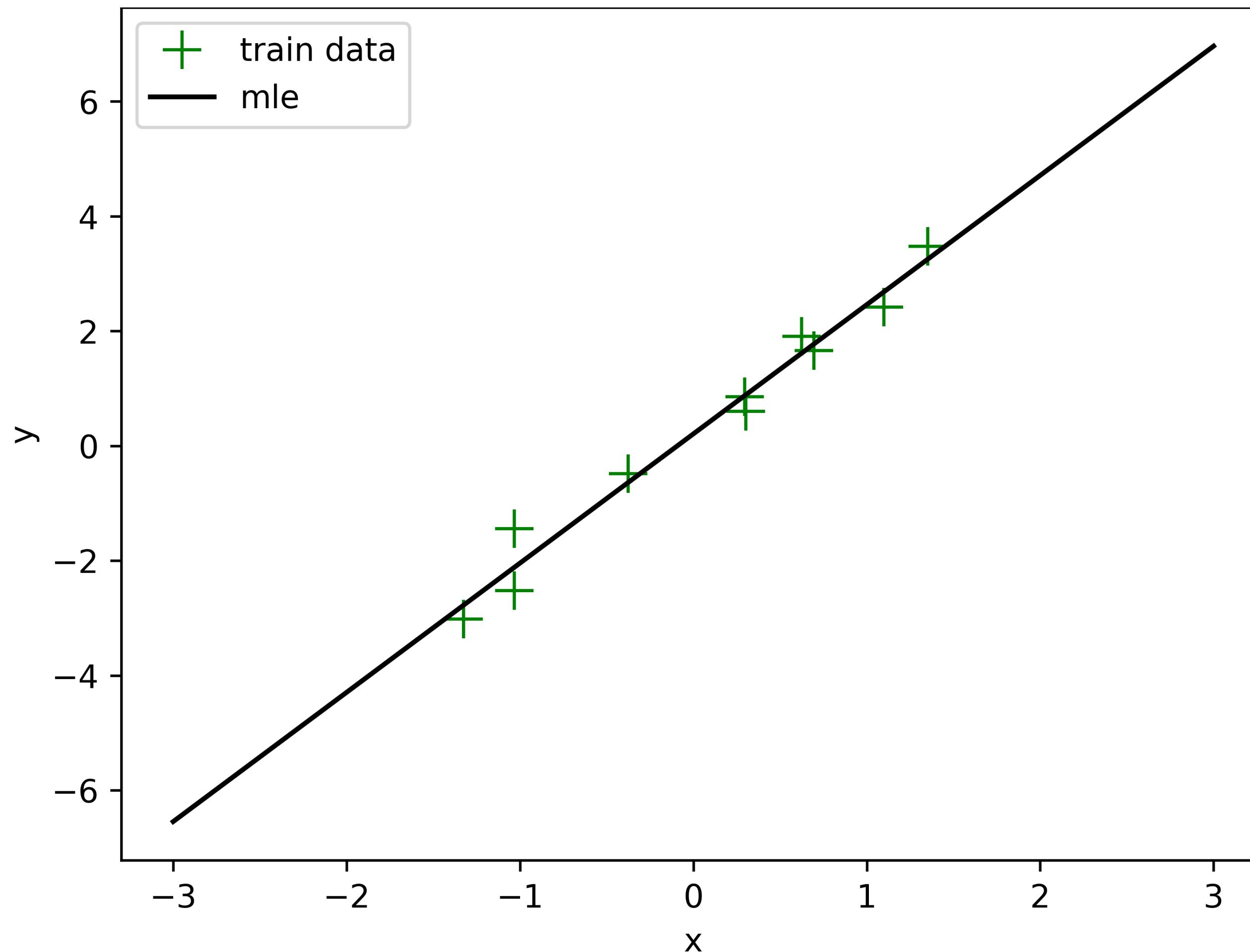
Mean = MAP

$$\Sigma = (\sigma^{-2} X^T X + \sigma_o^{-2} I_D)^{-1}$$

# Recap: Linear regression

Point estimate

$$\theta_{\text{ML}} = (X^T X)^{-1} X^T y \text{ or } \theta_{\text{MAP}} = (X^T X + N\lambda I)^{-1} X^T y$$



Exact posterior

$$p(\theta | \mathcal{D}) = \mathcal{N}(\theta; \mu, \Sigma)$$

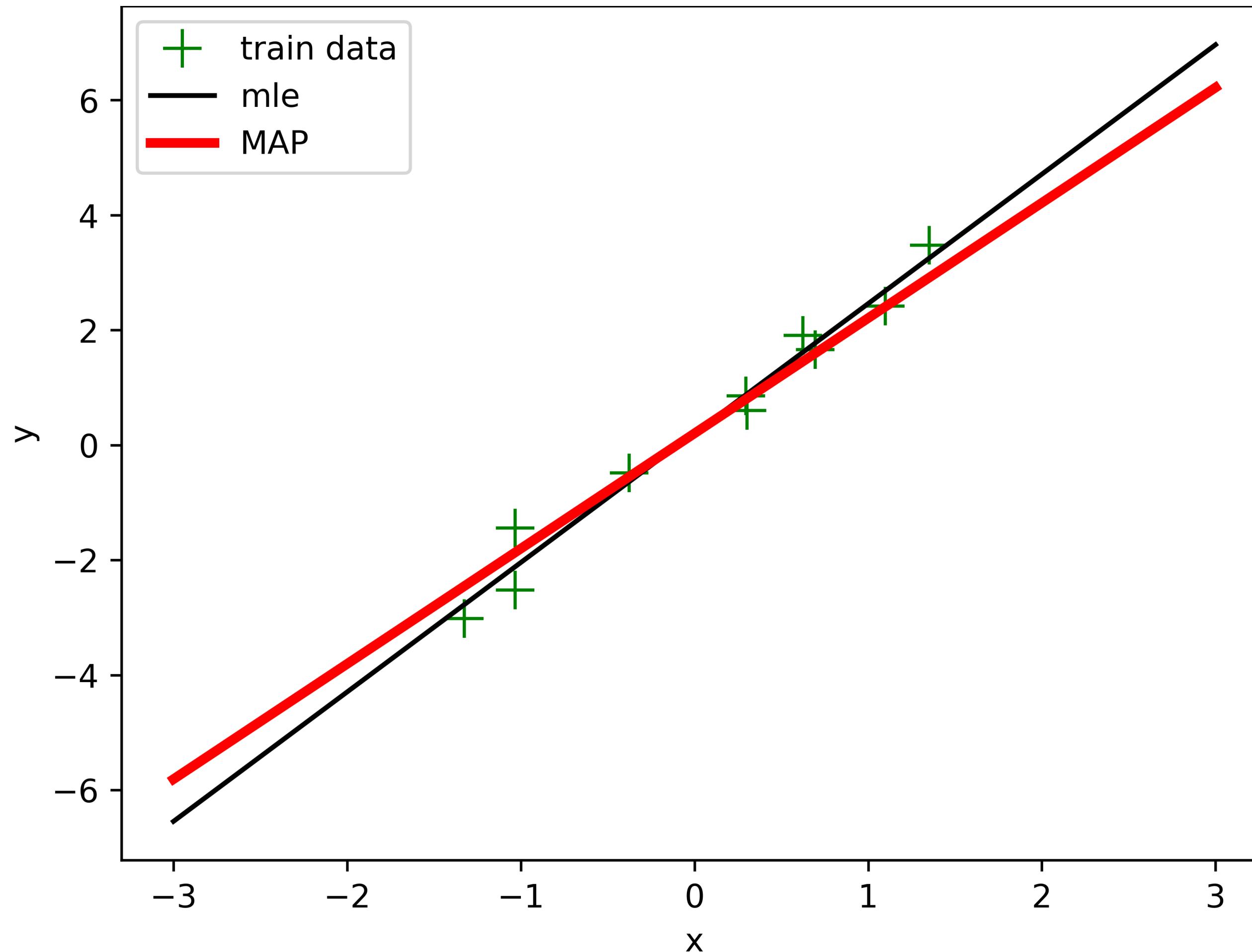
Mean = MAP

$$\mu = \Sigma \sigma^{-2} X^T y$$
$$\Sigma = (\sigma^{-2} X^T X + \sigma_o^{-2} I_D)^{-1}$$

# Recap: Linear regression

Point estimate

$$\theta_{\text{ML}} = (X^T X)^{-1} X^T \mathbf{y} \text{ or } \theta_{\text{MAP}} = (X^T X + N\lambda \mathbf{I})^{-1} X^T \mathbf{y}$$



Exact posterior

$$p(\theta | \mathcal{D}) = \mathcal{N}(\theta; \mu, \Sigma)$$

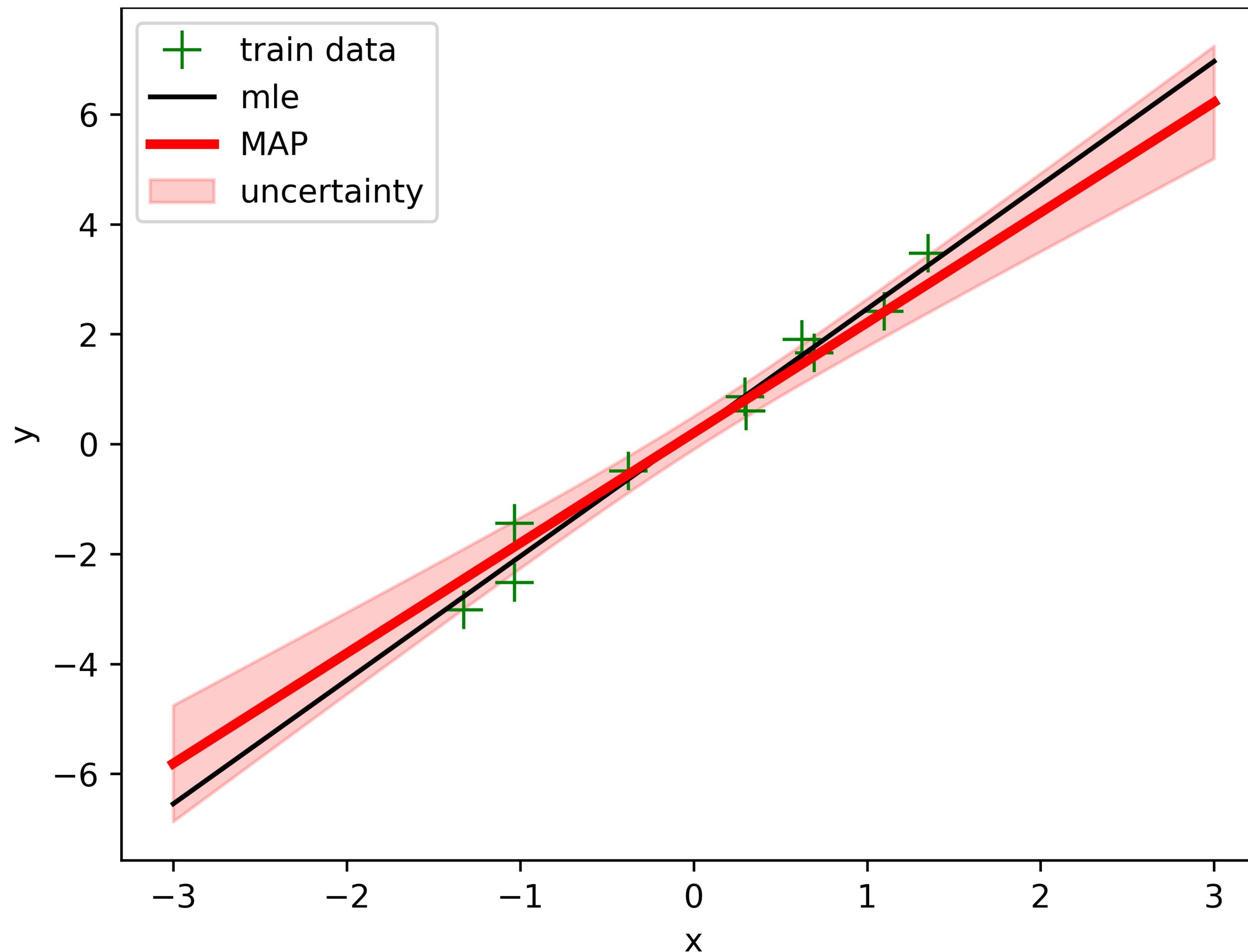
Mean = MAP

$$\mu = \Sigma \sigma^{-2} X^T \mathbf{y}$$
$$\Sigma = (\sigma^{-2} X^T X + \sigma_o^{-2} I_D)^{-1}$$

# Recap: Linear regression

Point estimate

$$\theta_{\text{ML}} = (X^T X)^{-1} X^T y \text{ or } \theta_{\text{MAP}} = (X^T X + N\lambda I)^{-1} X^T y$$



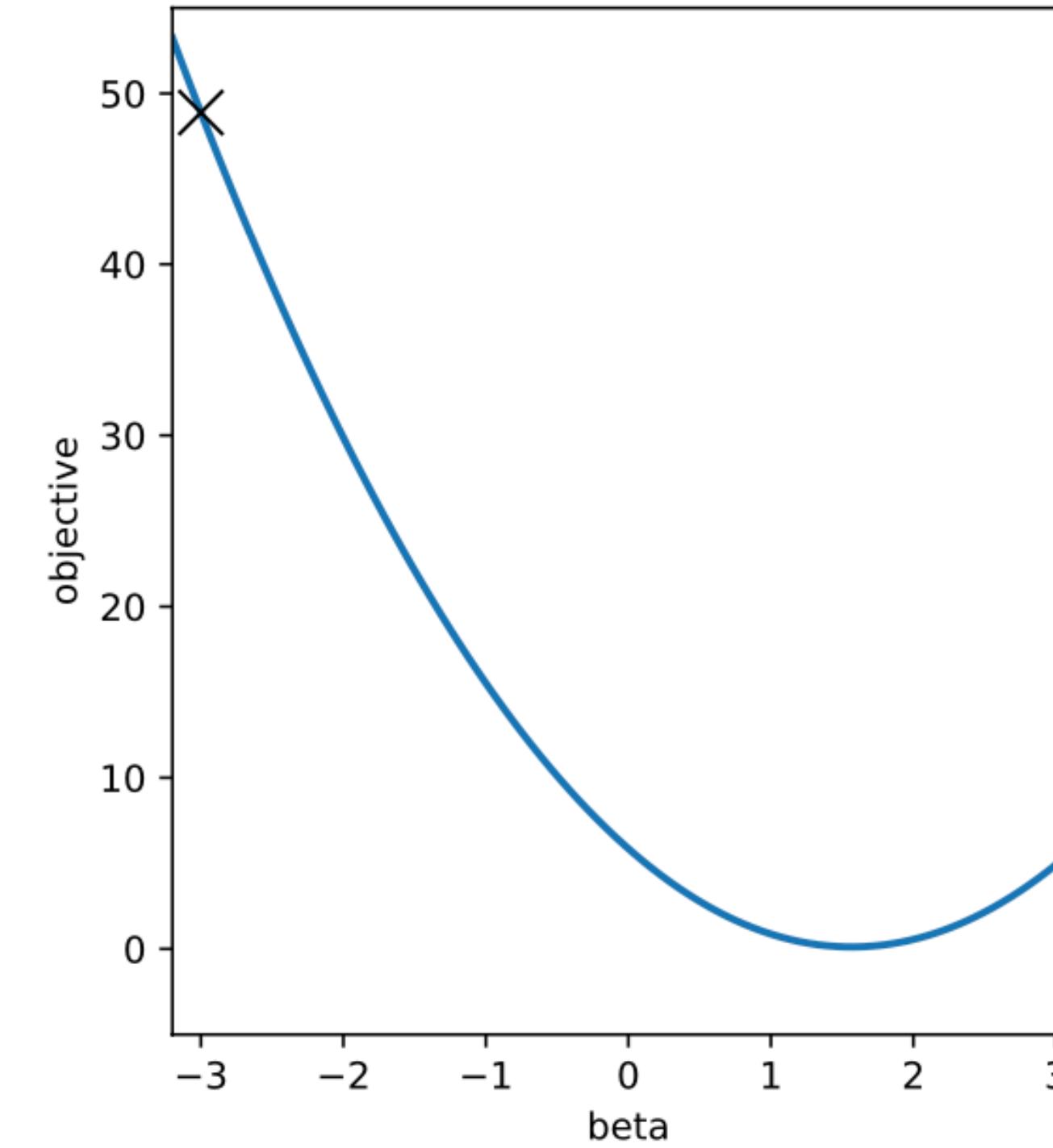
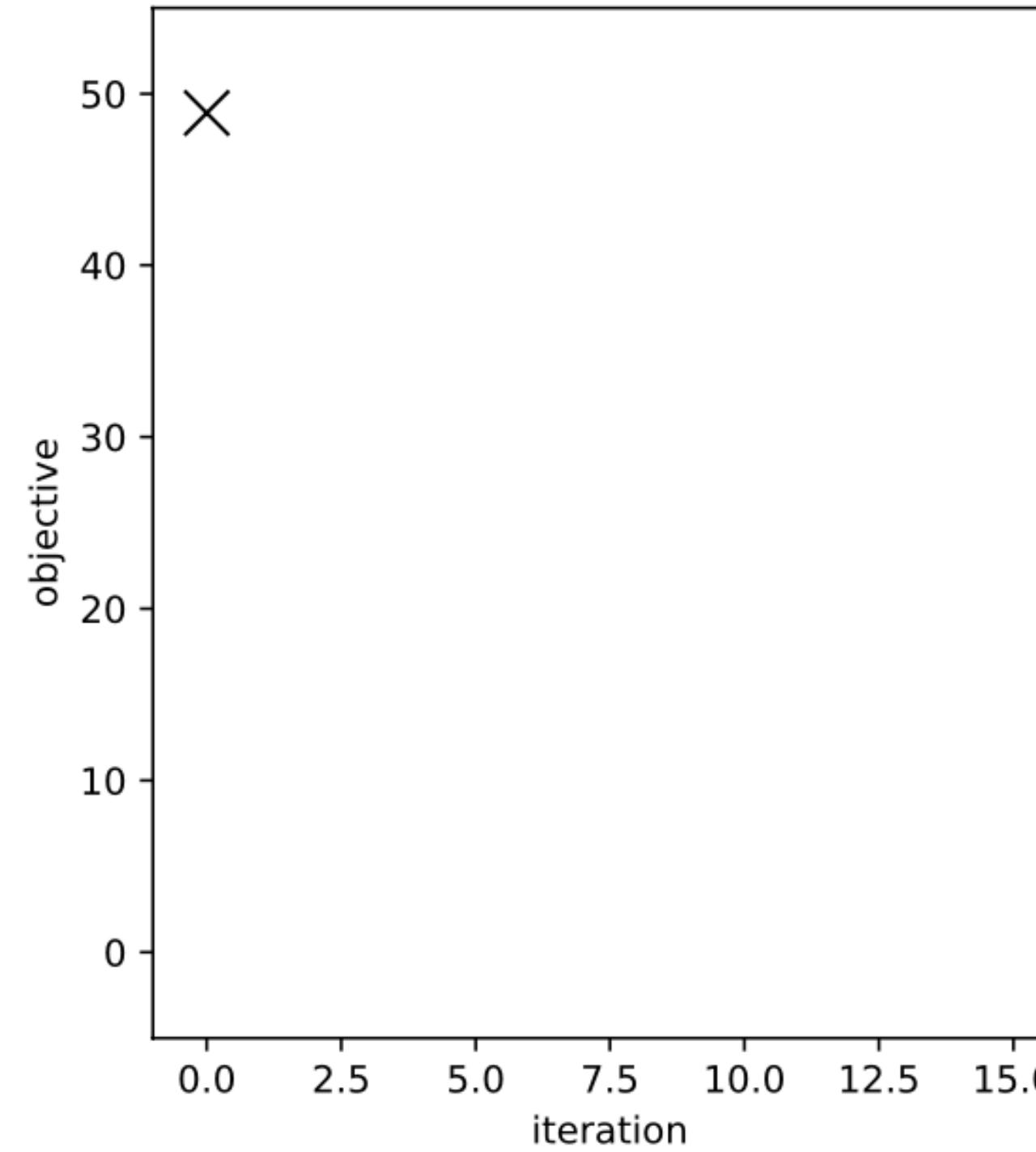
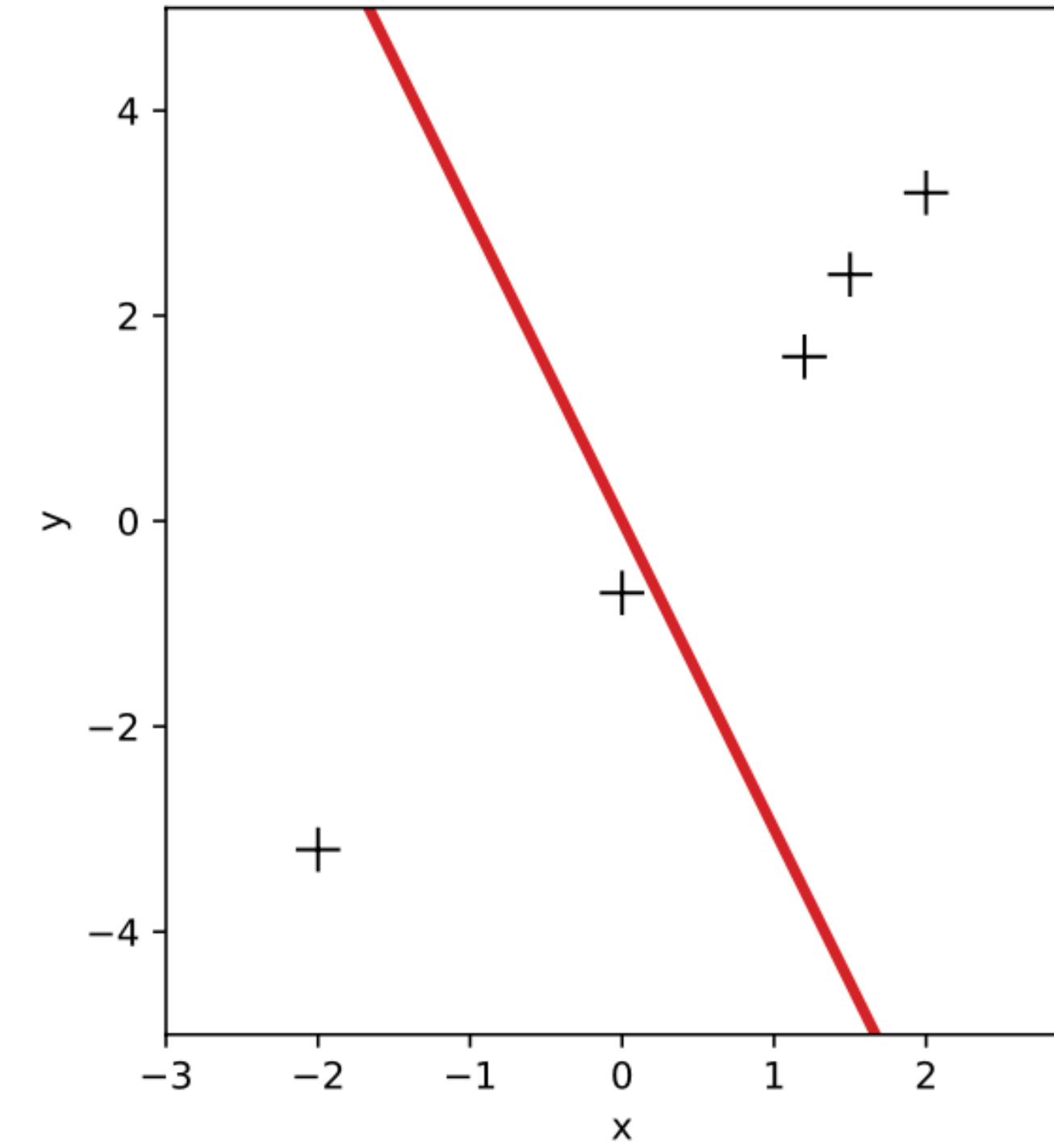
Exact posterior

$$p(\theta | \mathcal{D}) = \mathcal{N}(\theta; \mu, \Sigma)$$

Mean = MAP

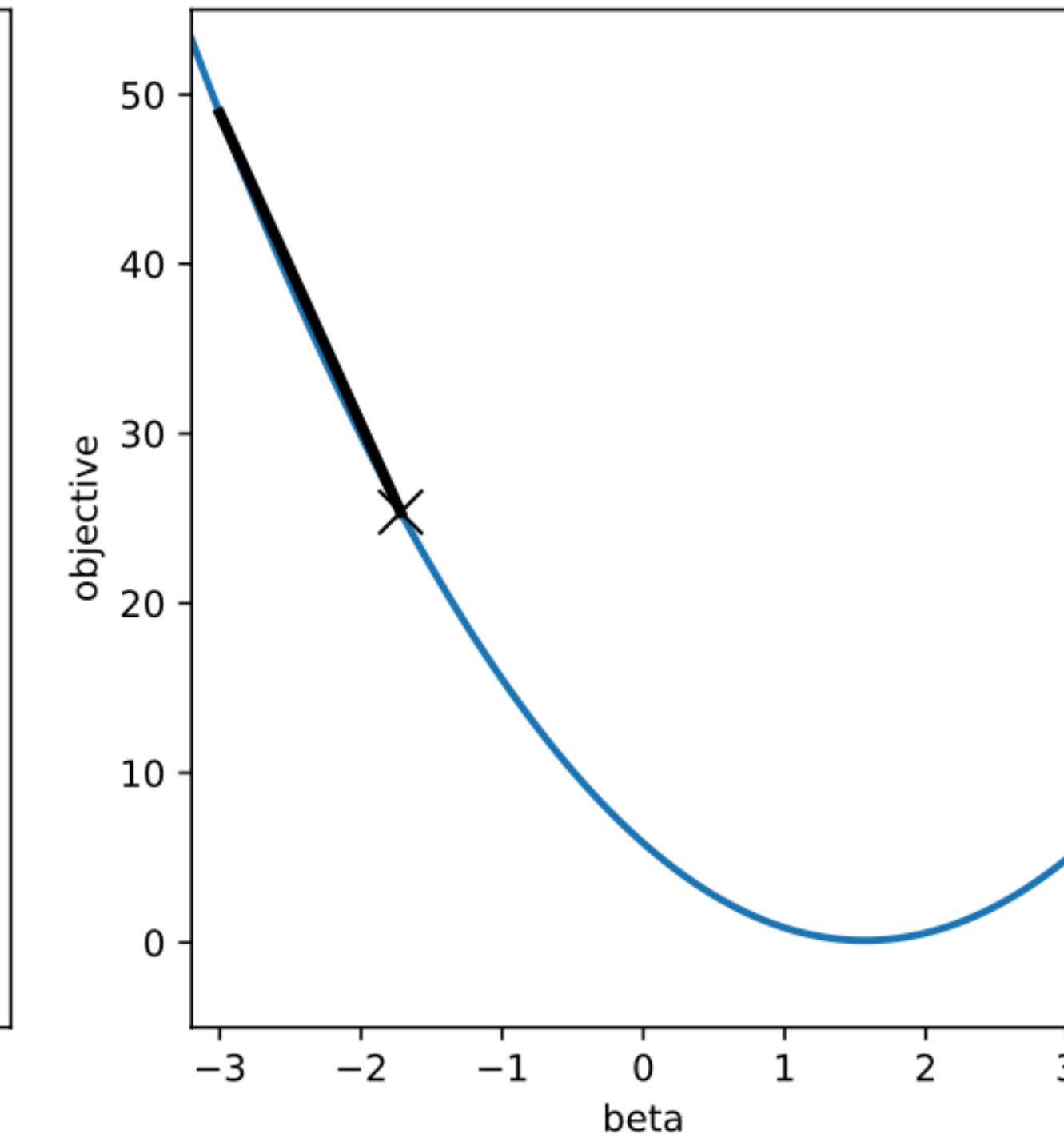
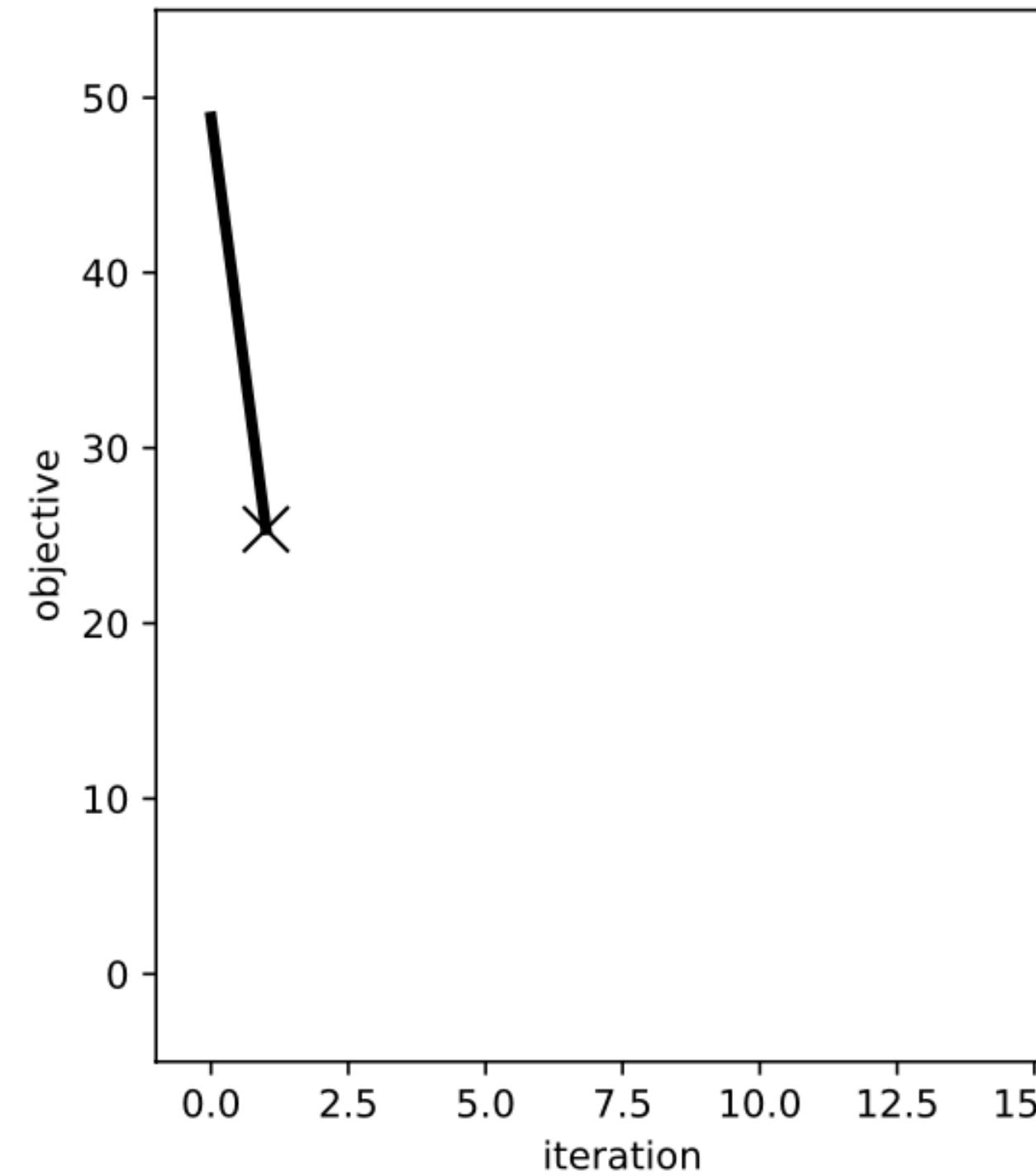
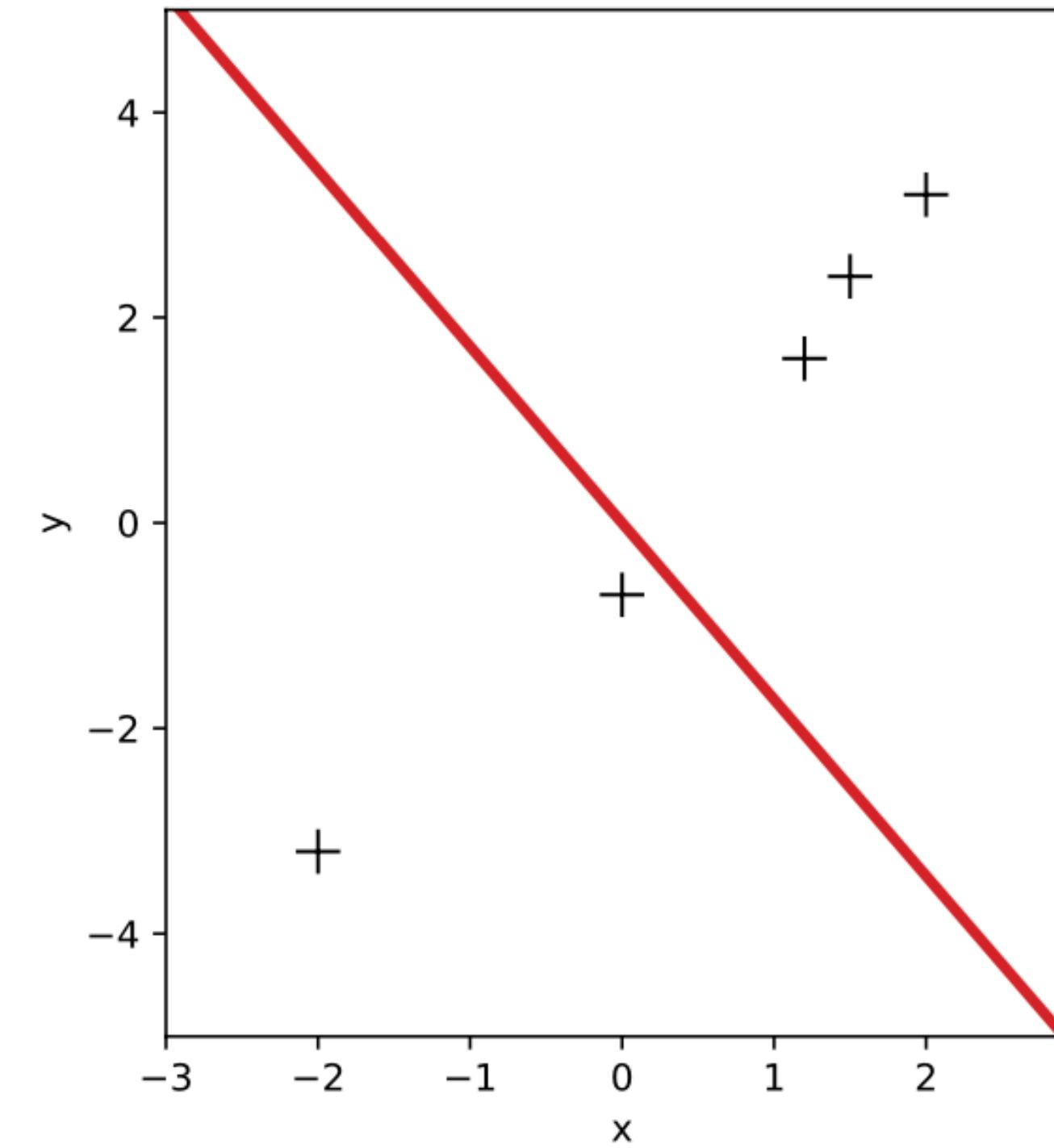
$$\mu = \Sigma \sigma^{-2} X^T y$$
$$\Sigma = (\sigma^{-2} X^T X + \sigma_o^{-2} I_D)^{-1}$$

# Recap: Optimisation for linear regression



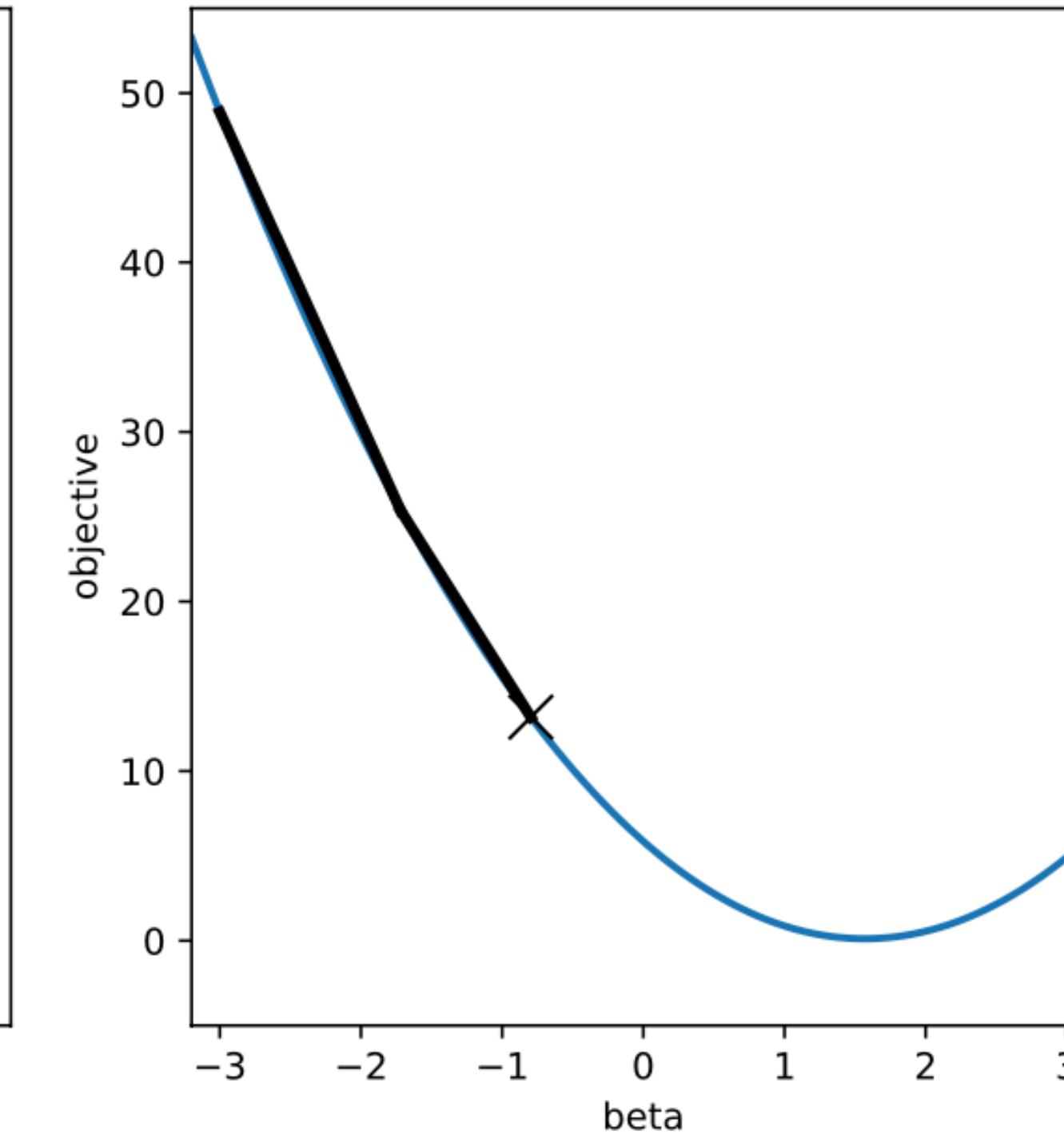
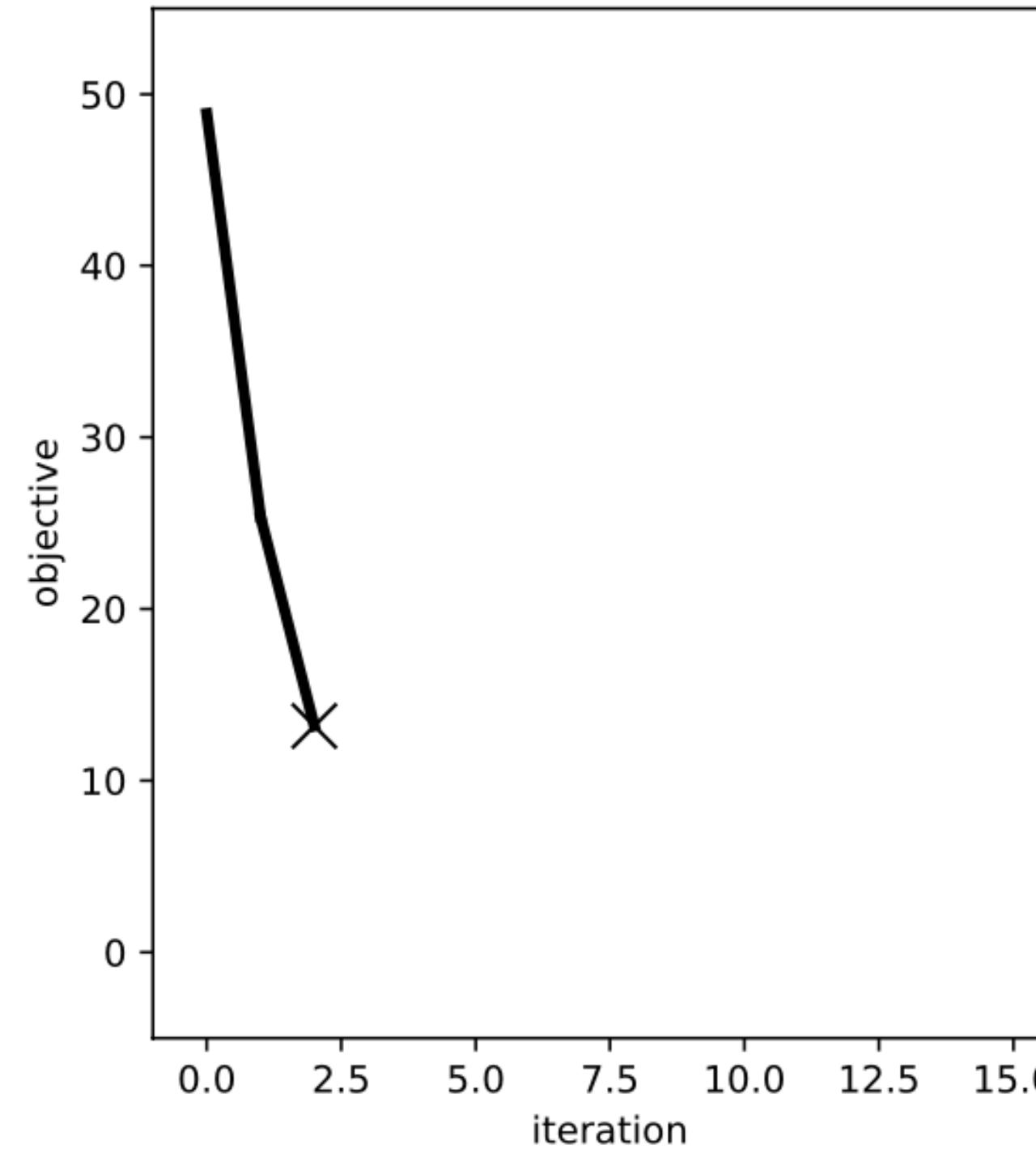
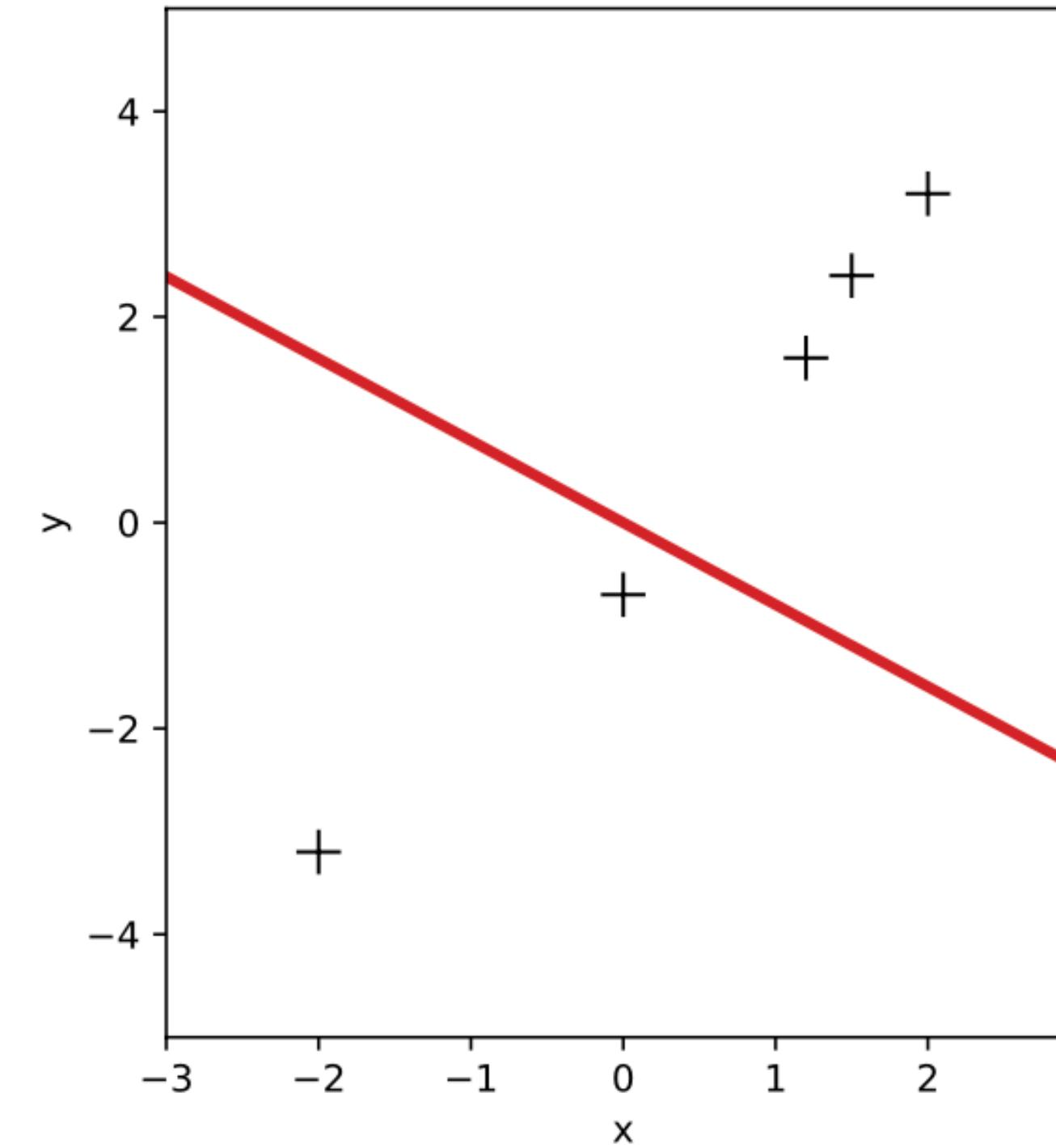
Another example: <https://playground.tensorflow.org/>

# Recap: Optimisation for linear regression



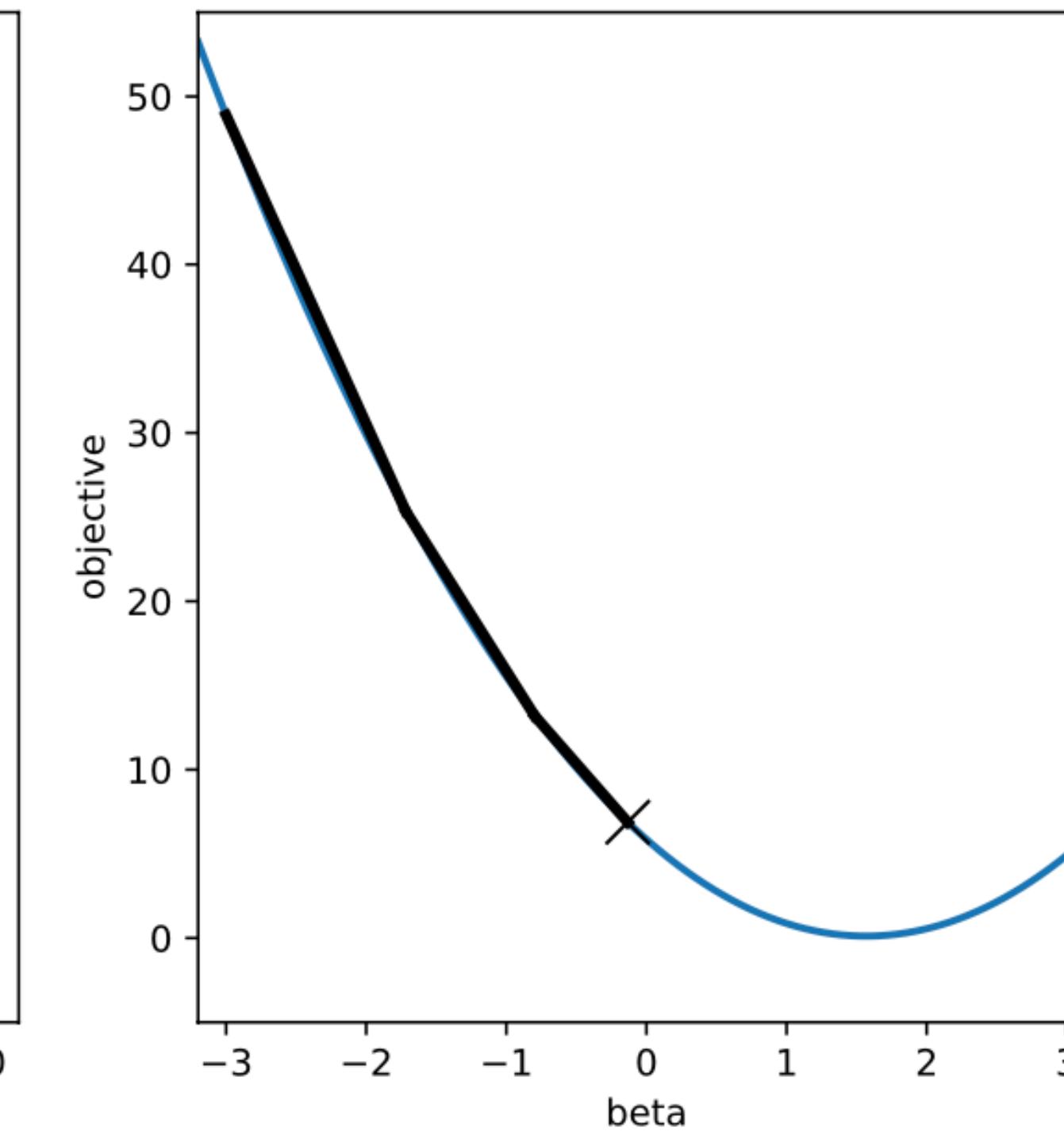
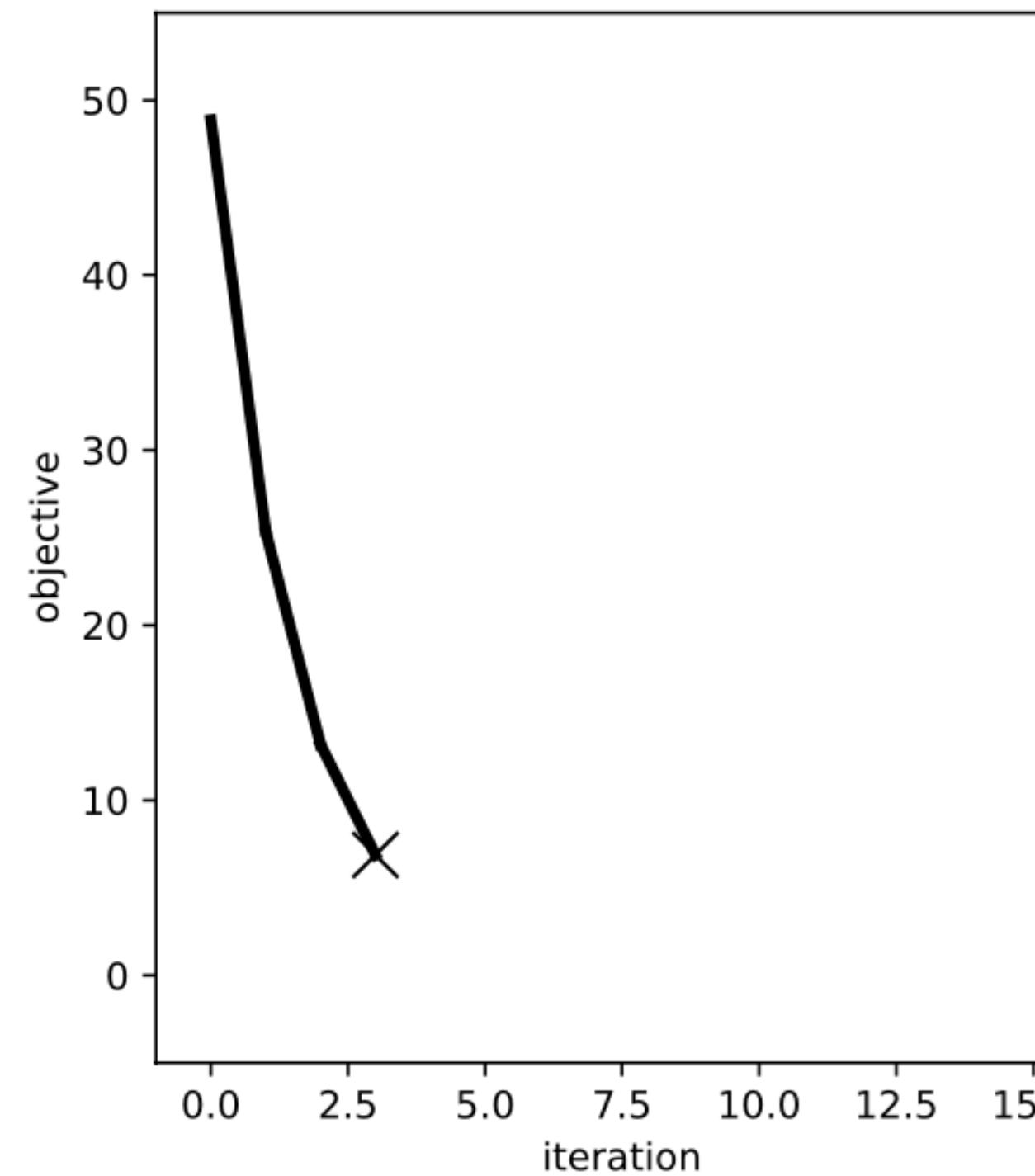
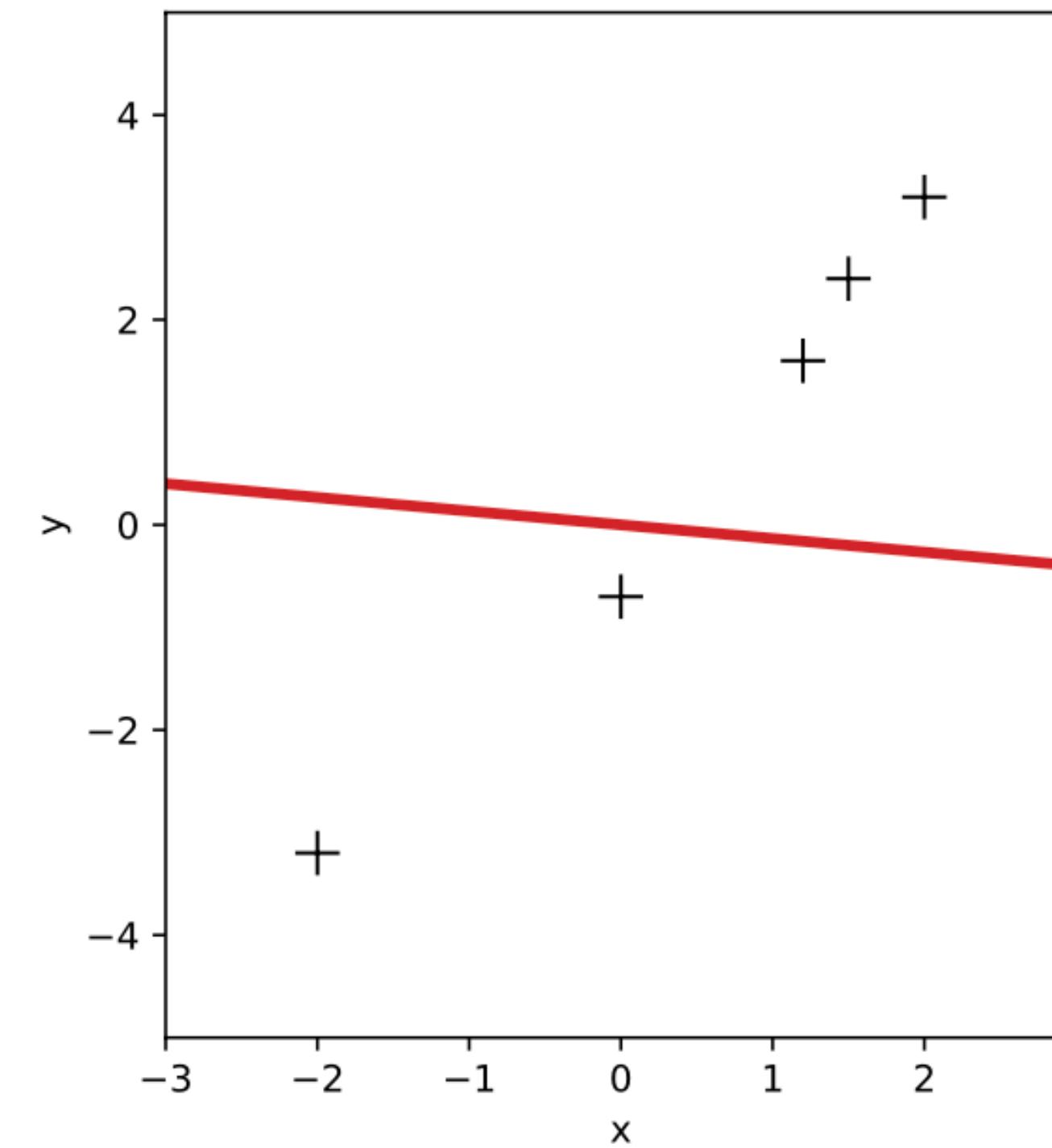
Another example: <https://playground.tensorflow.org/>

# Recap: Optimisation for linear regression



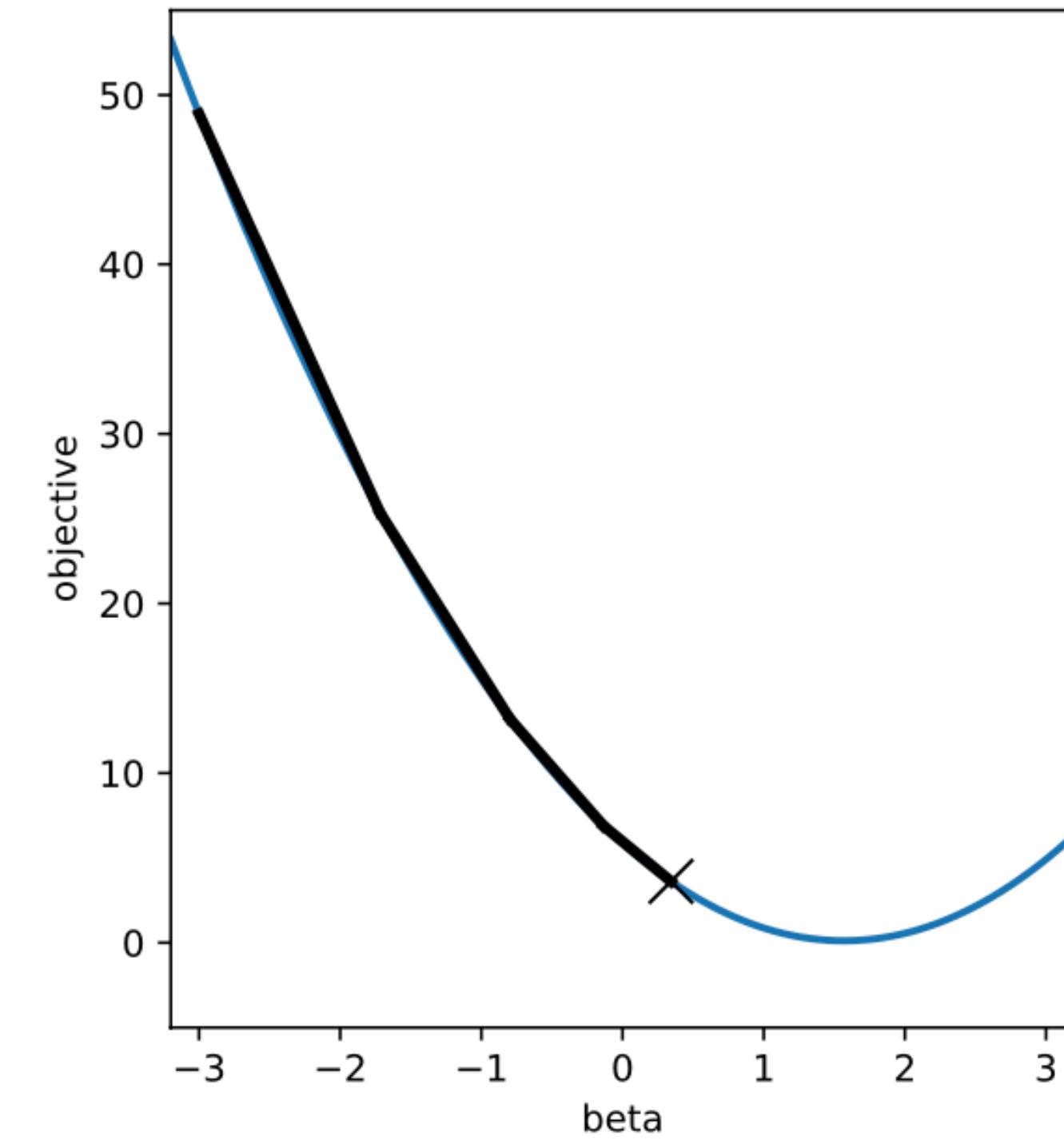
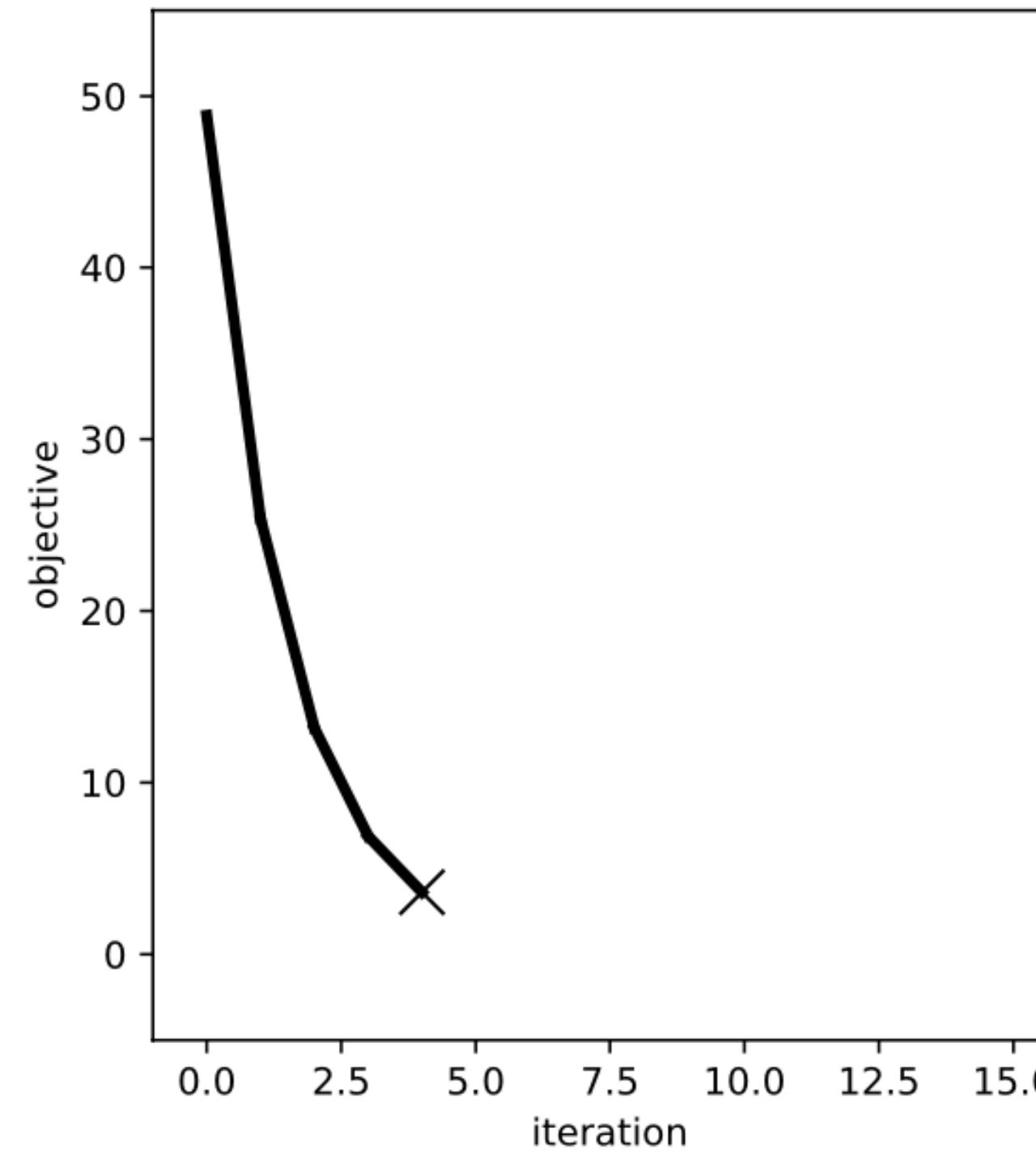
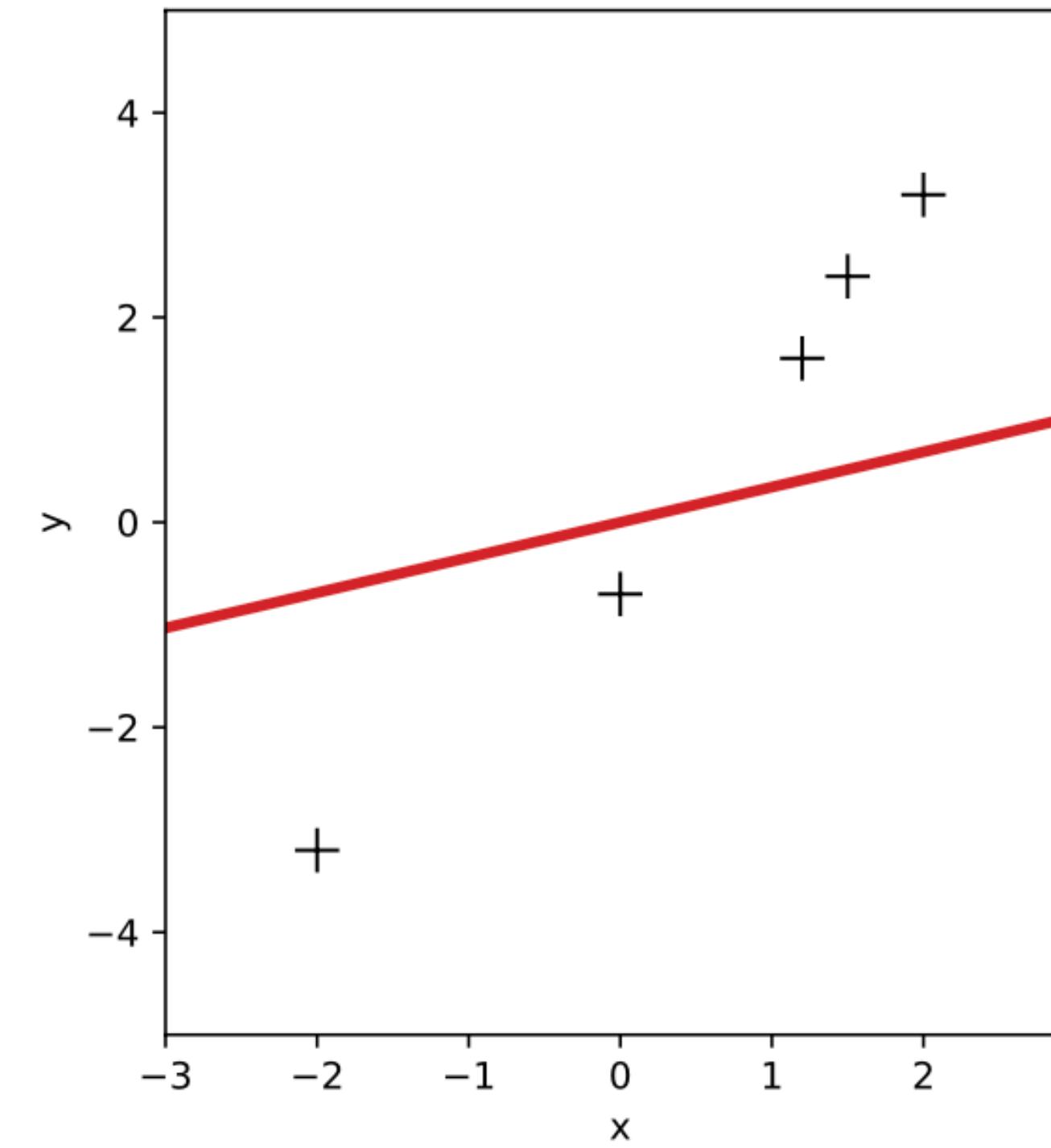
Another example: <https://playground.tensorflow.org/>

# Recap: Optimisation for linear regression



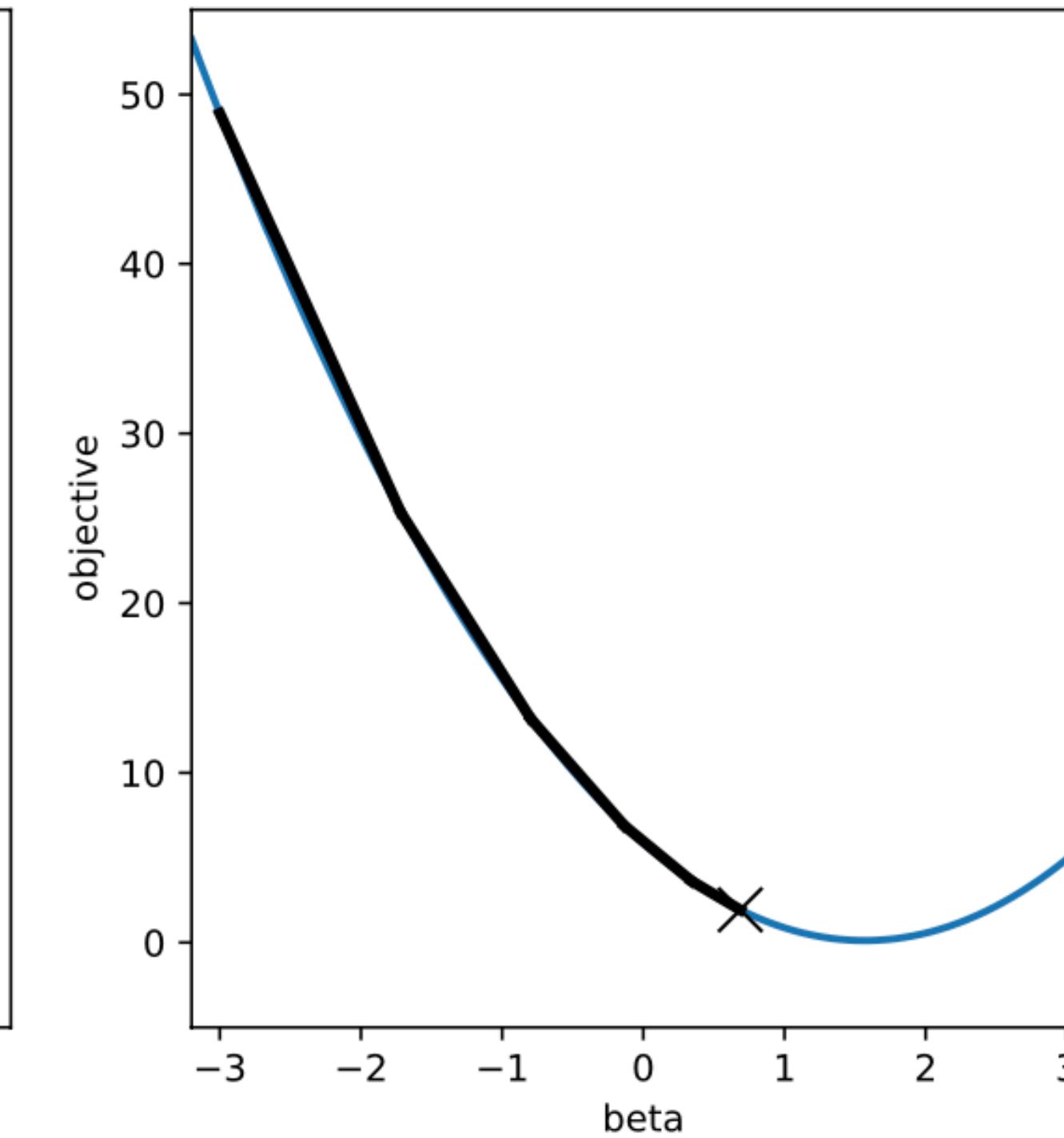
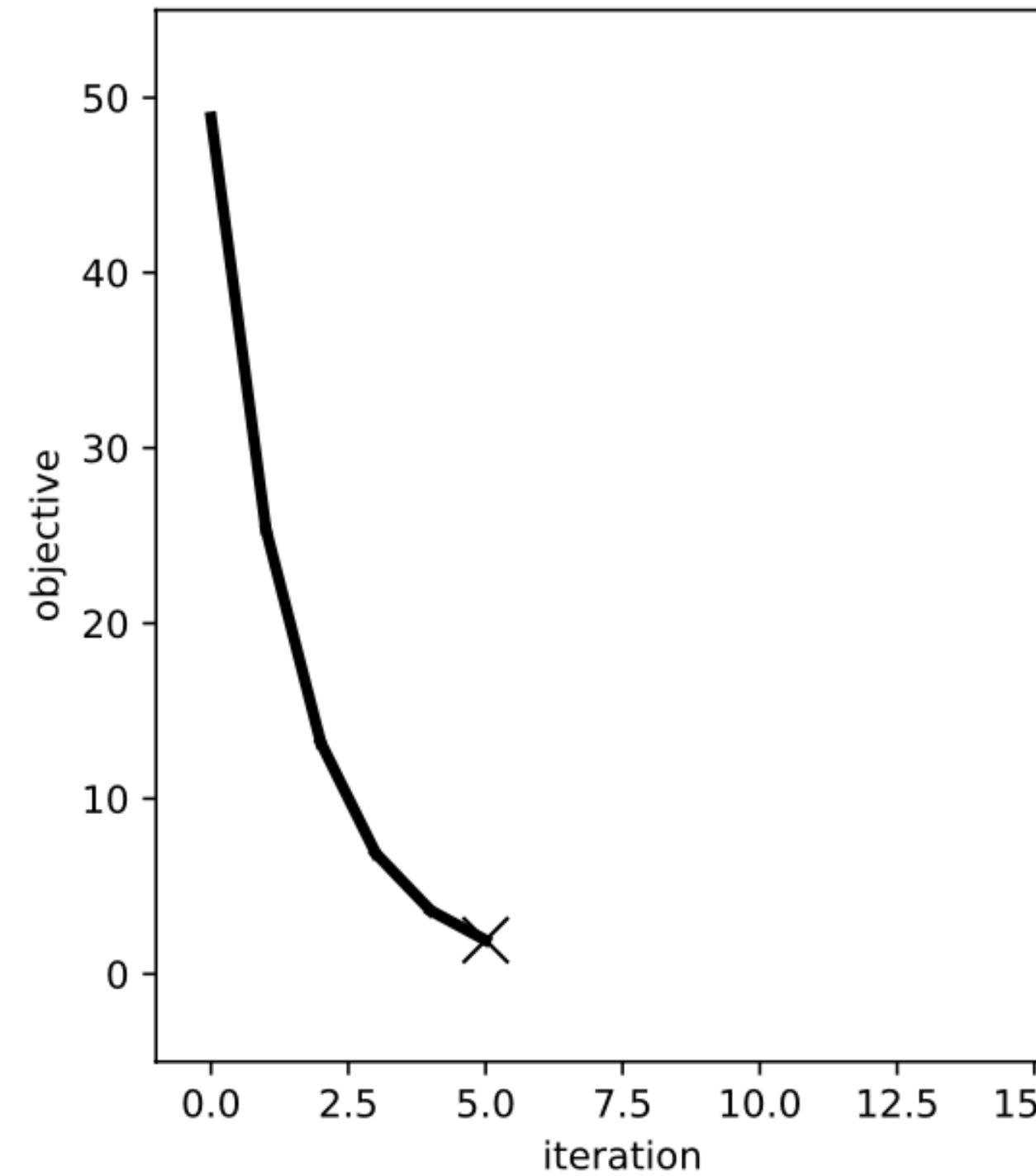
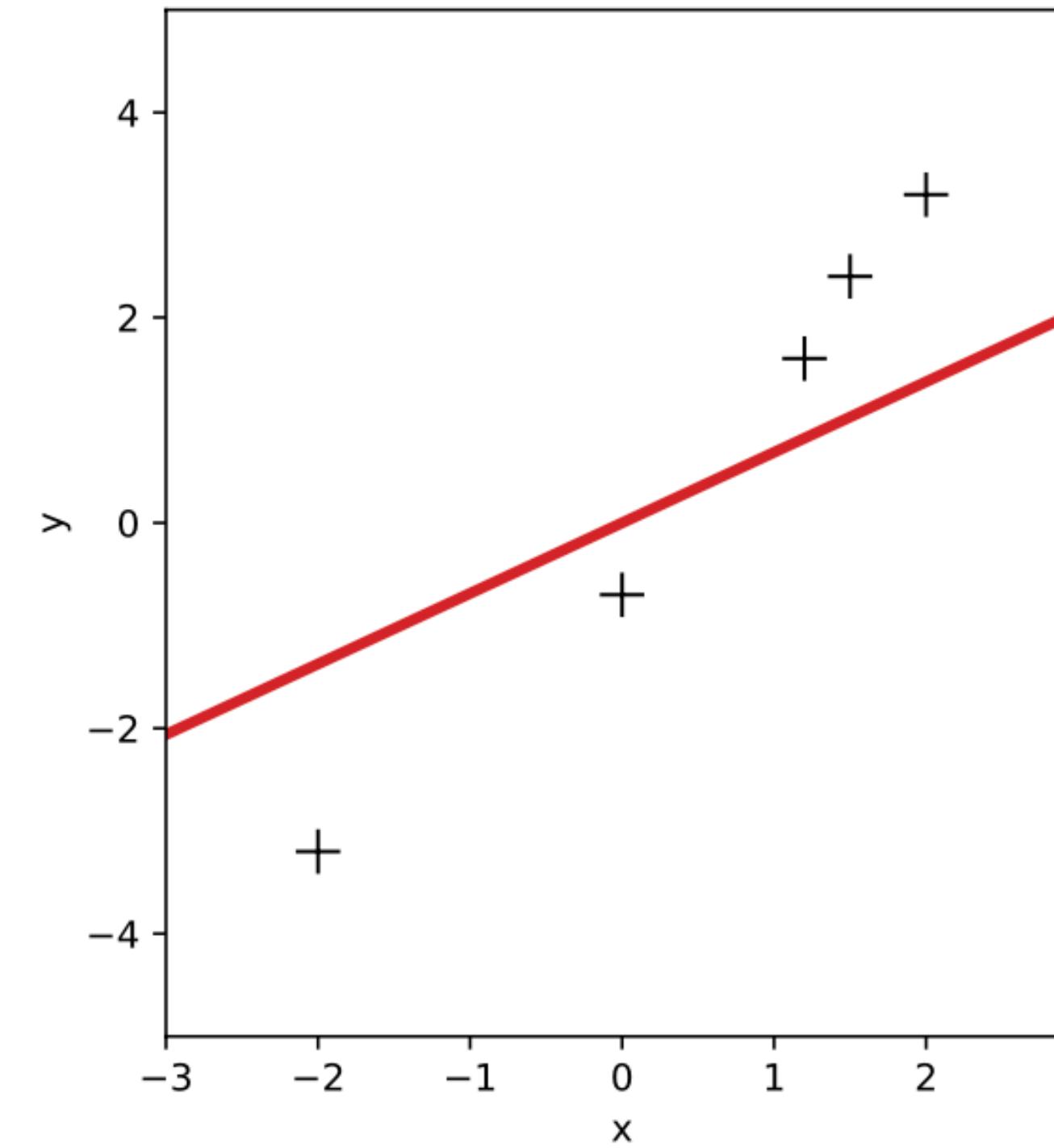
Another example: <https://playground.tensorflow.org/>

# Recap: Optimisation for linear regression



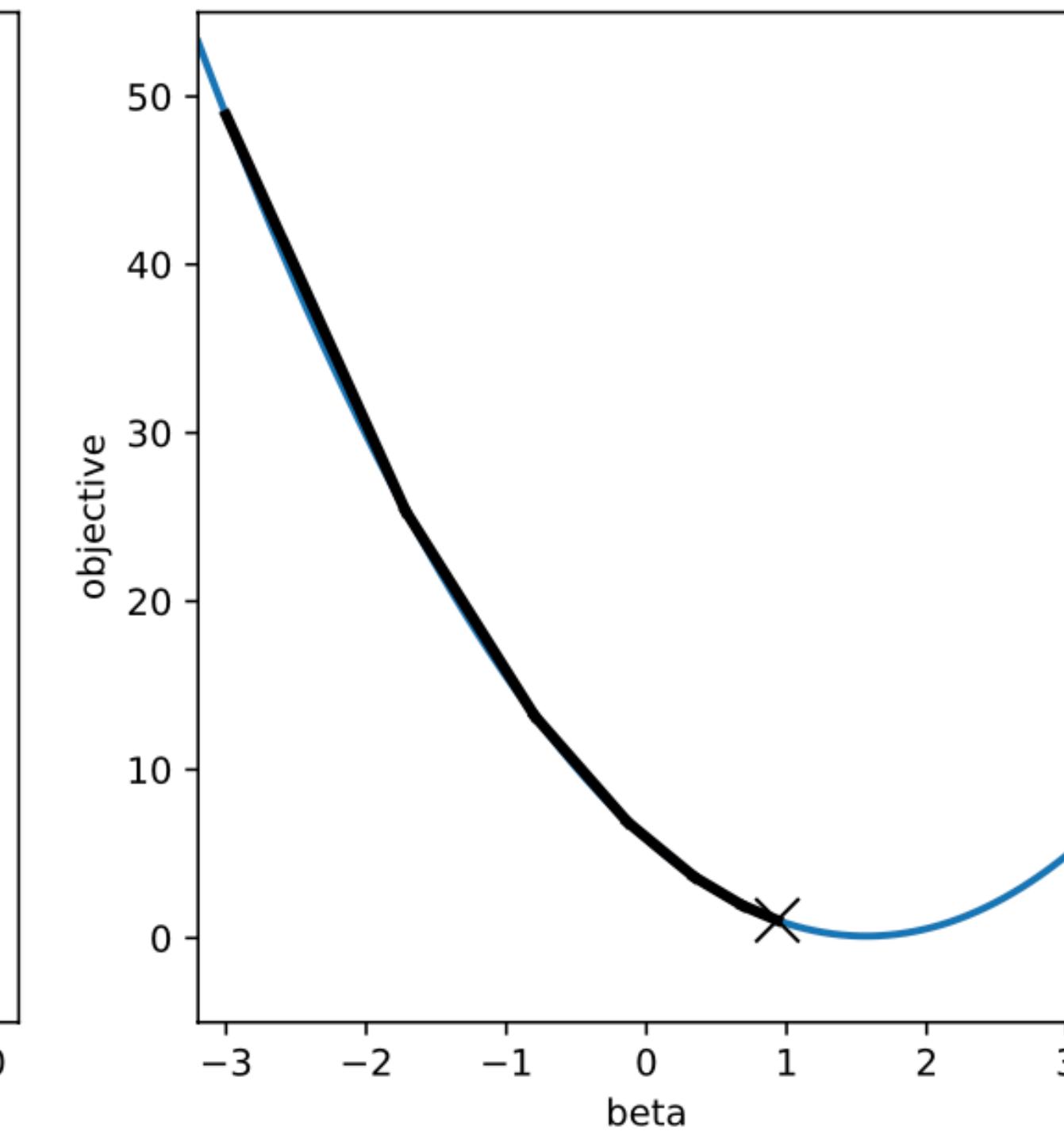
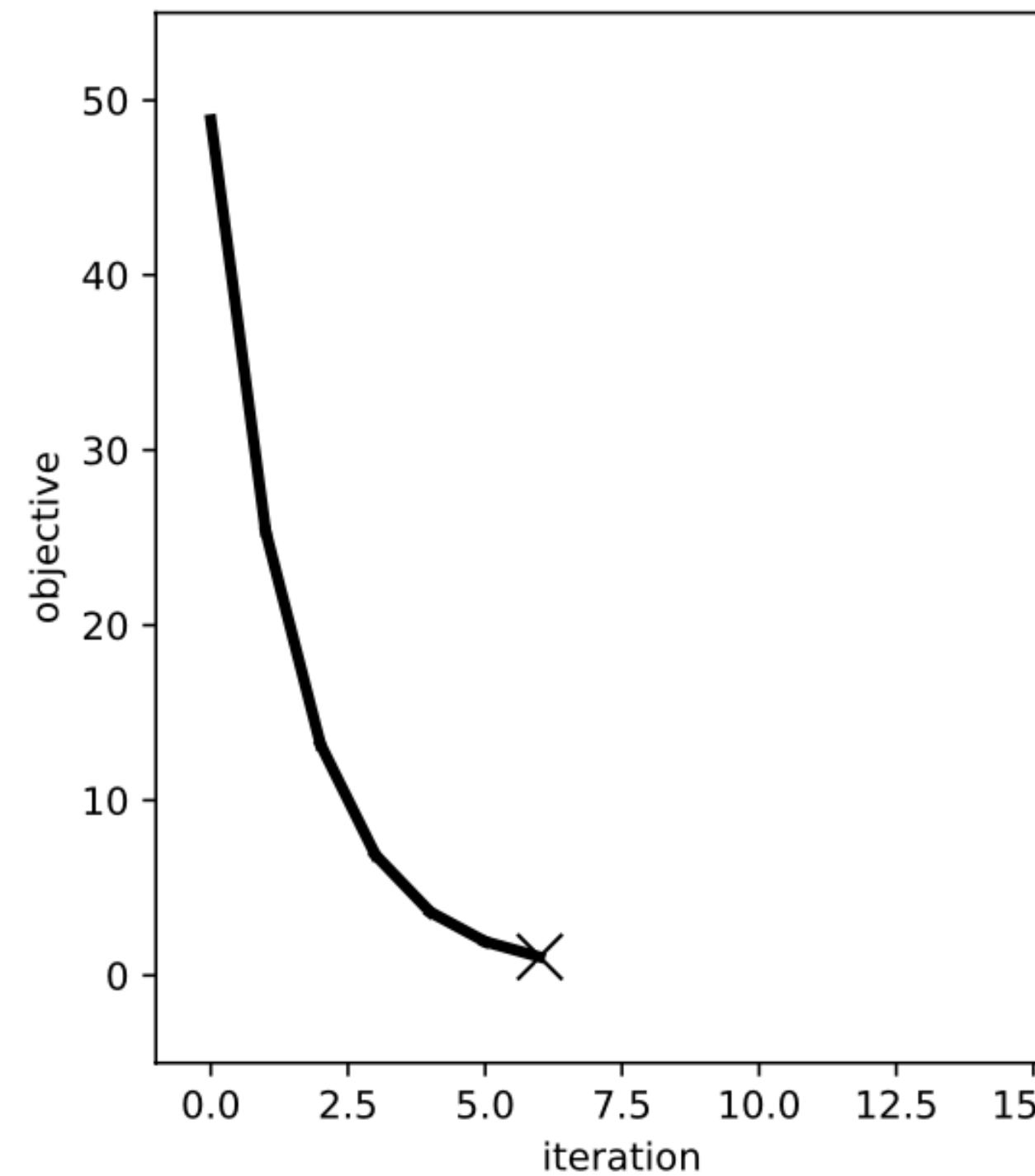
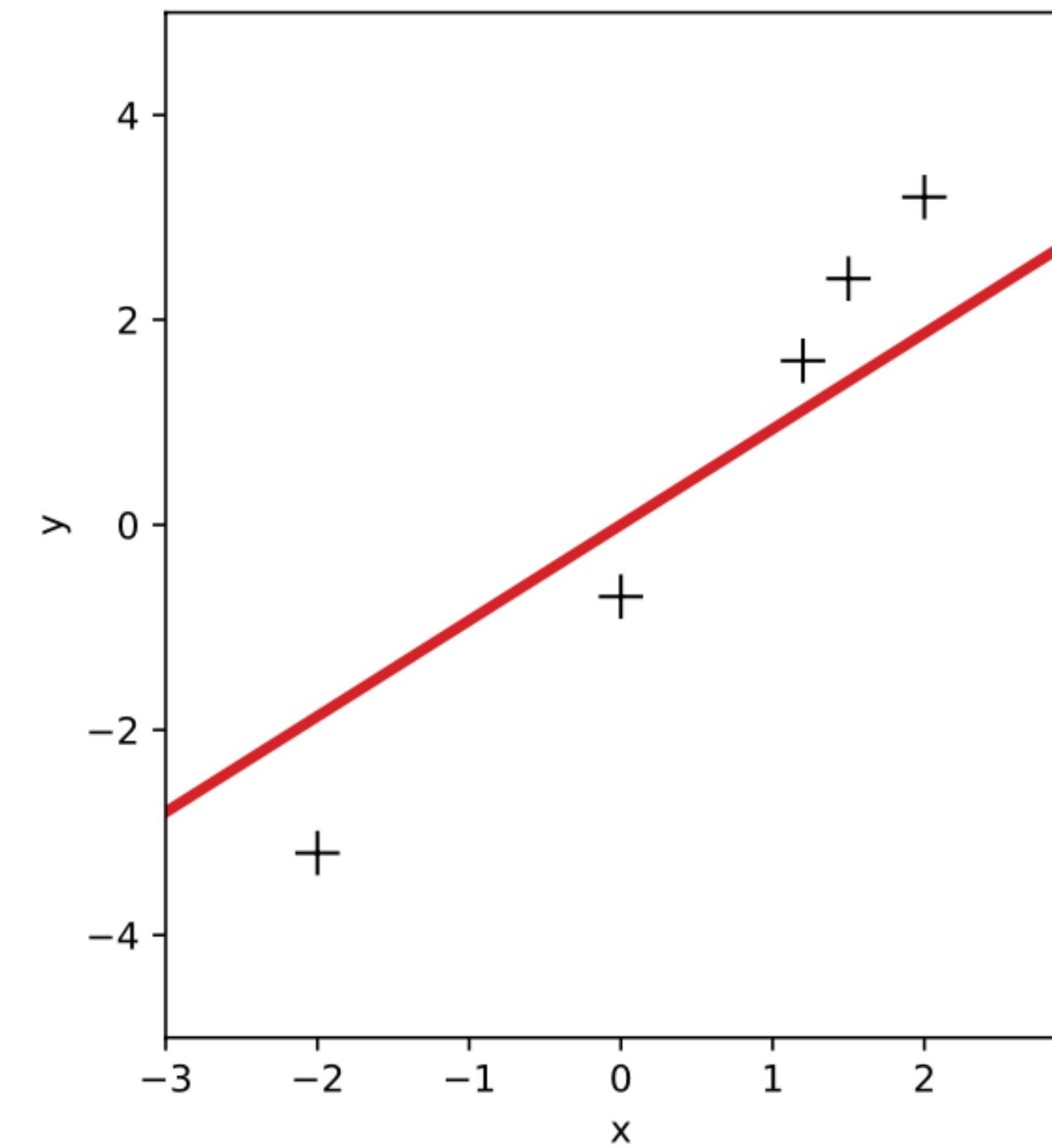
Another example: <https://playground.tensorflow.org/>

# Recap: Optimisation for linear regression



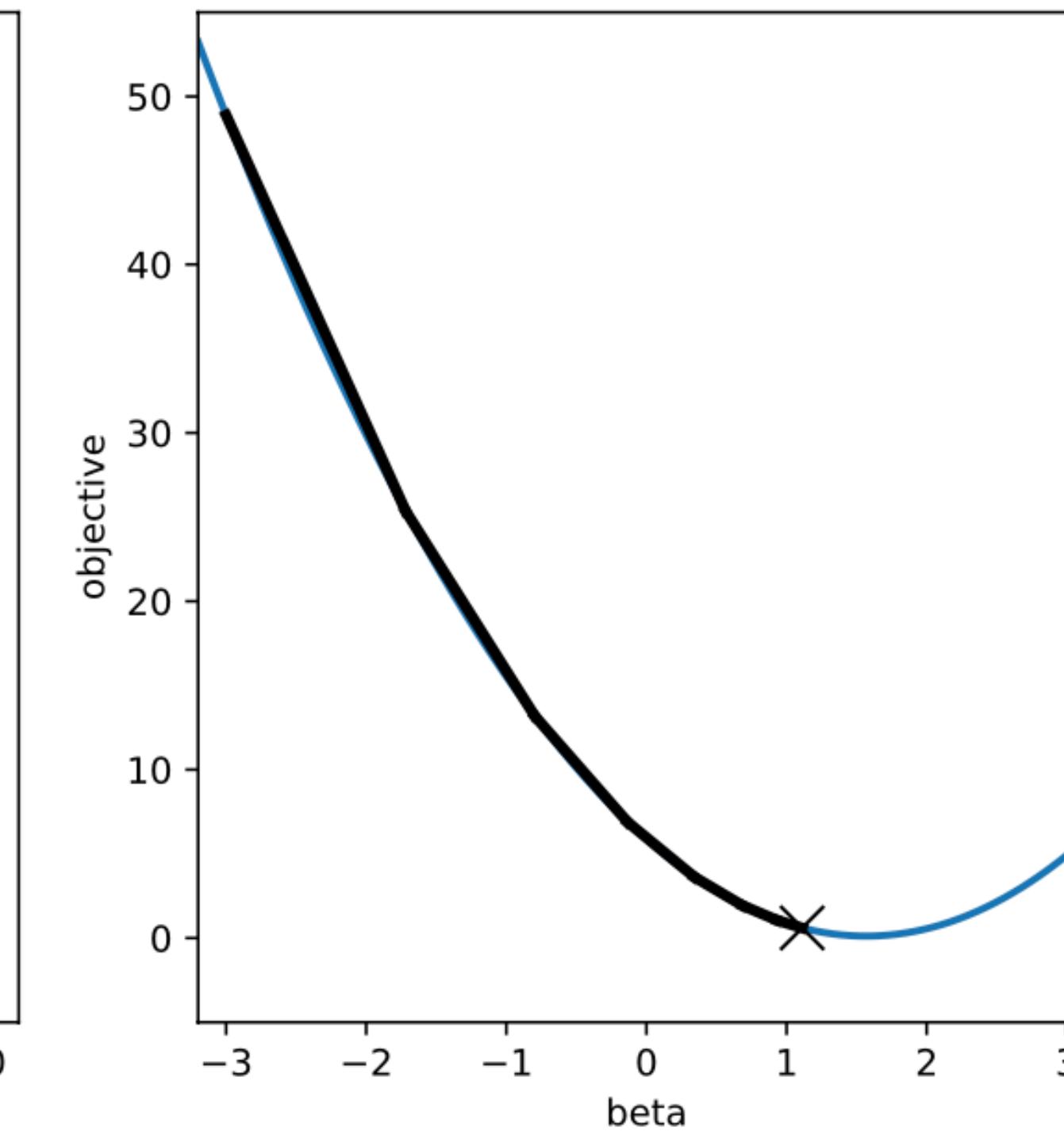
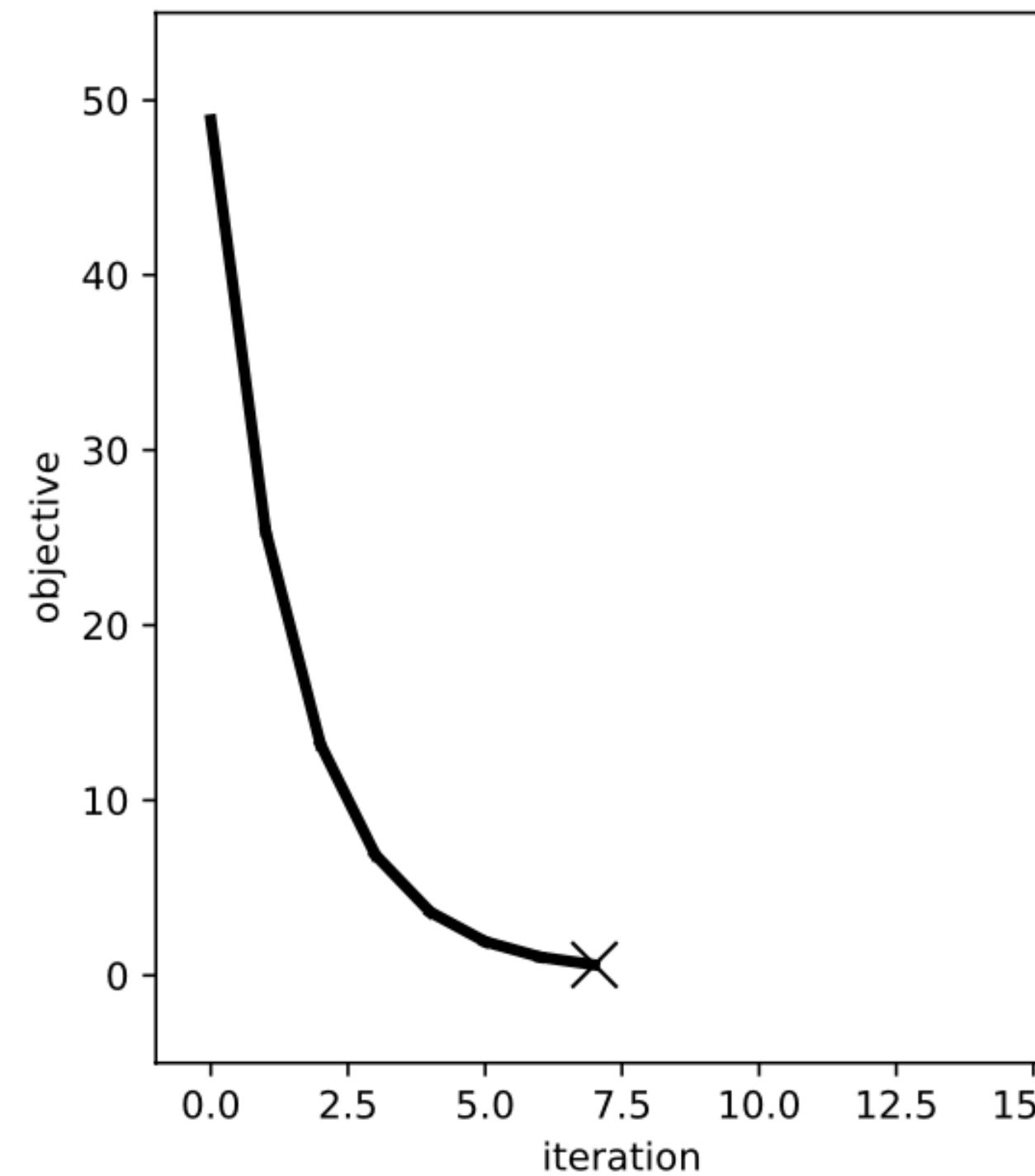
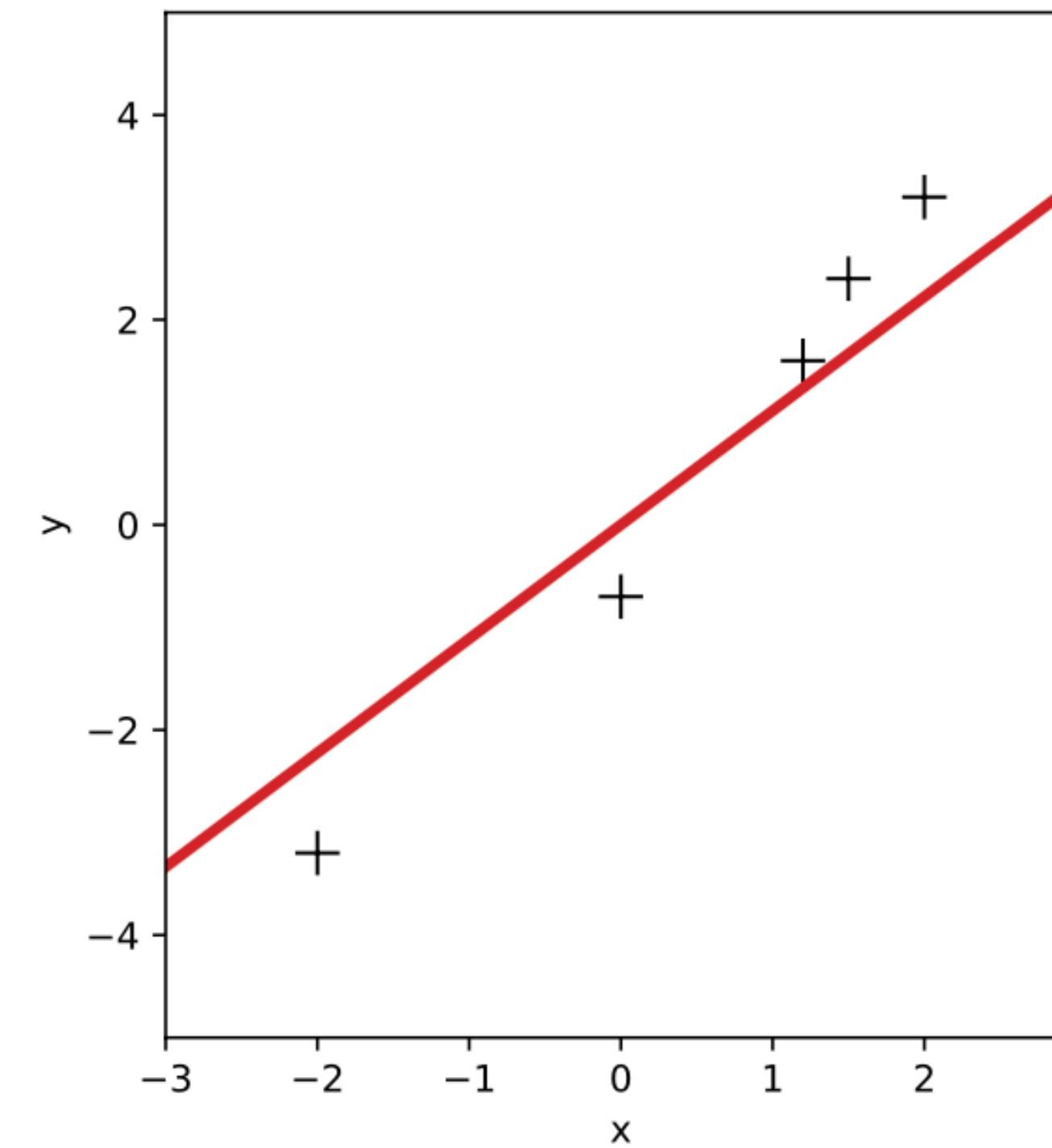
Another example: <https://playground.tensorflow.org/>

# Recap: Optimisation for linear regression



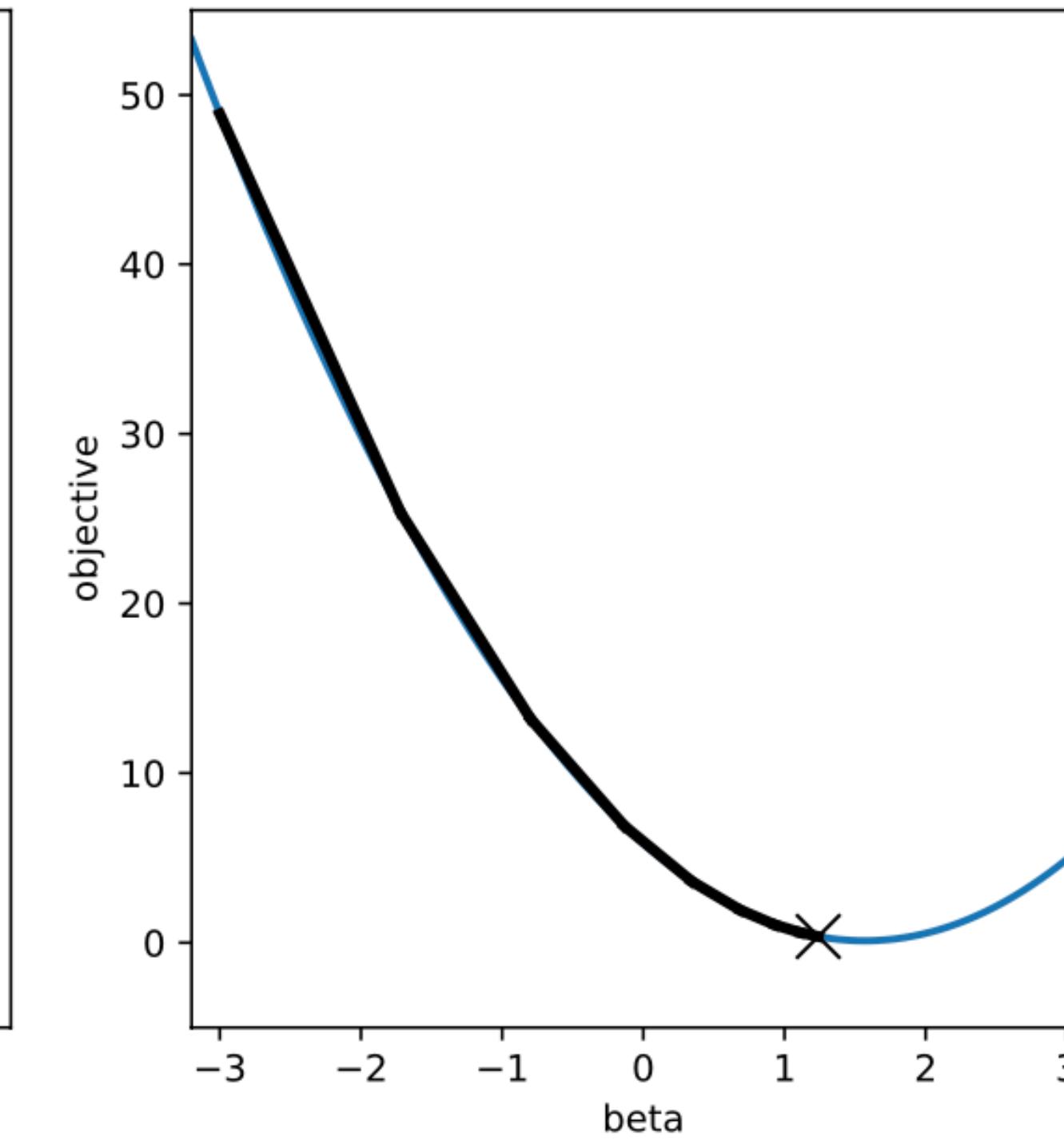
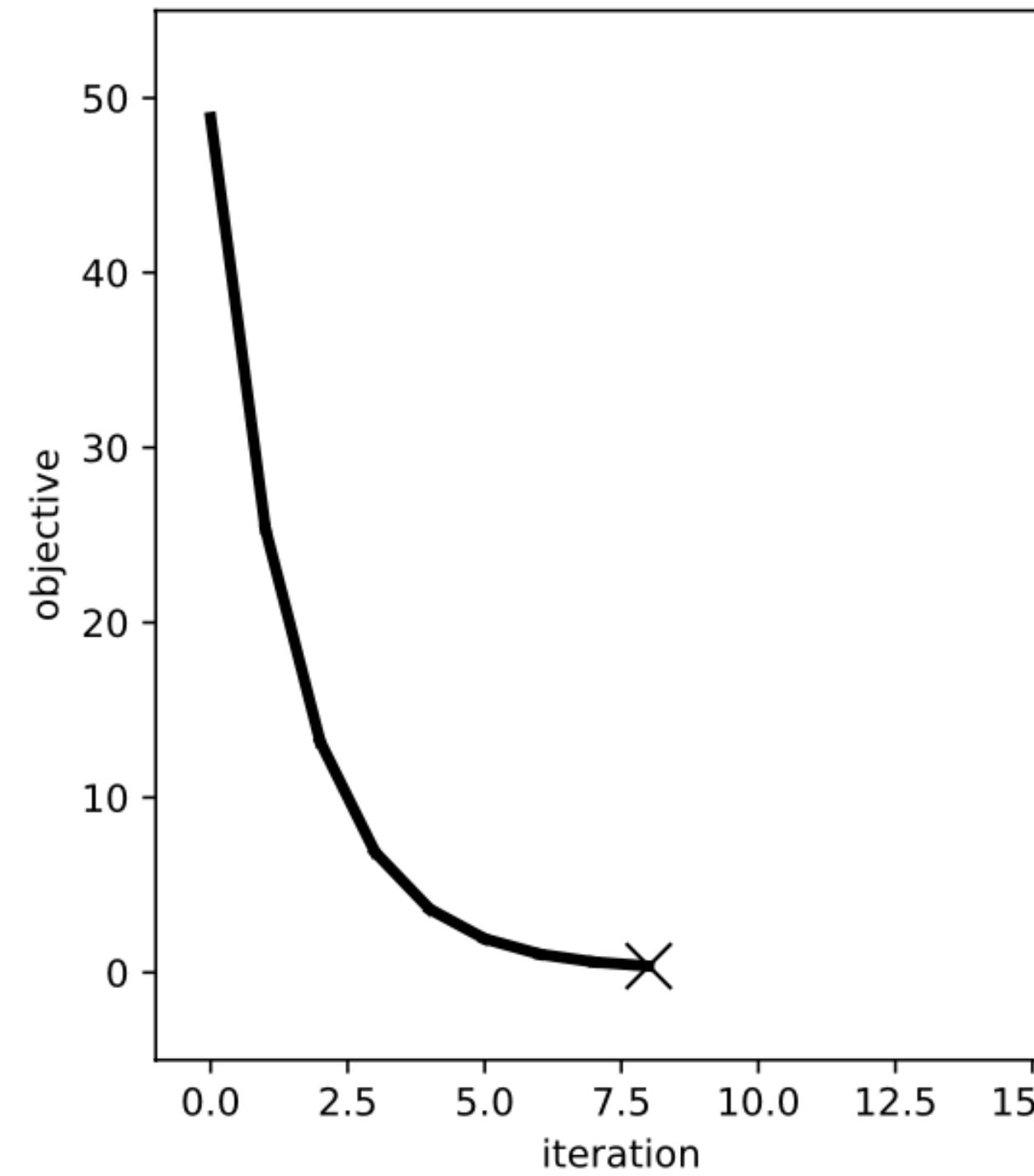
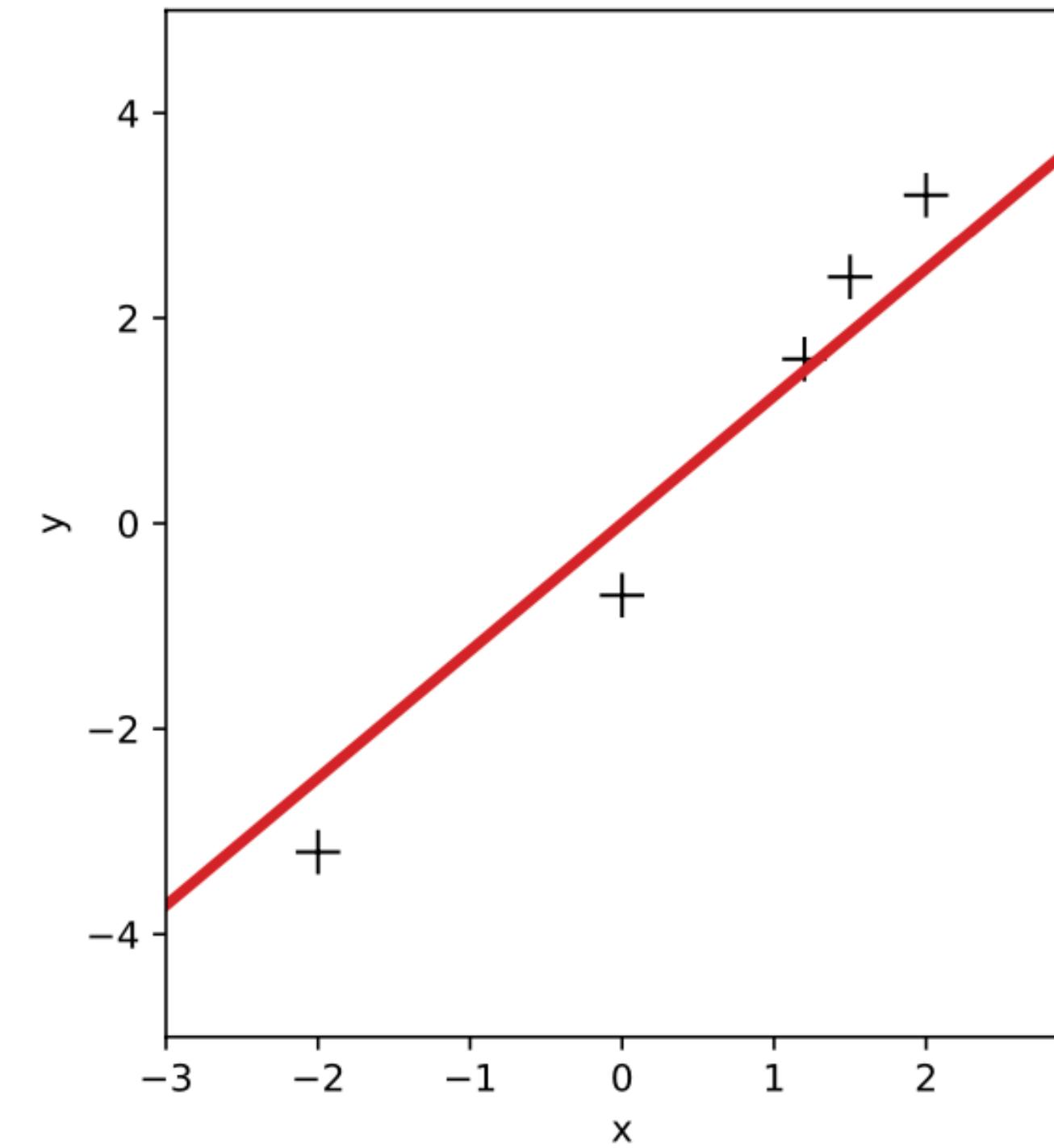
Another example: <https://playground.tensorflow.org/>

# Recap: Optimisation for linear regression



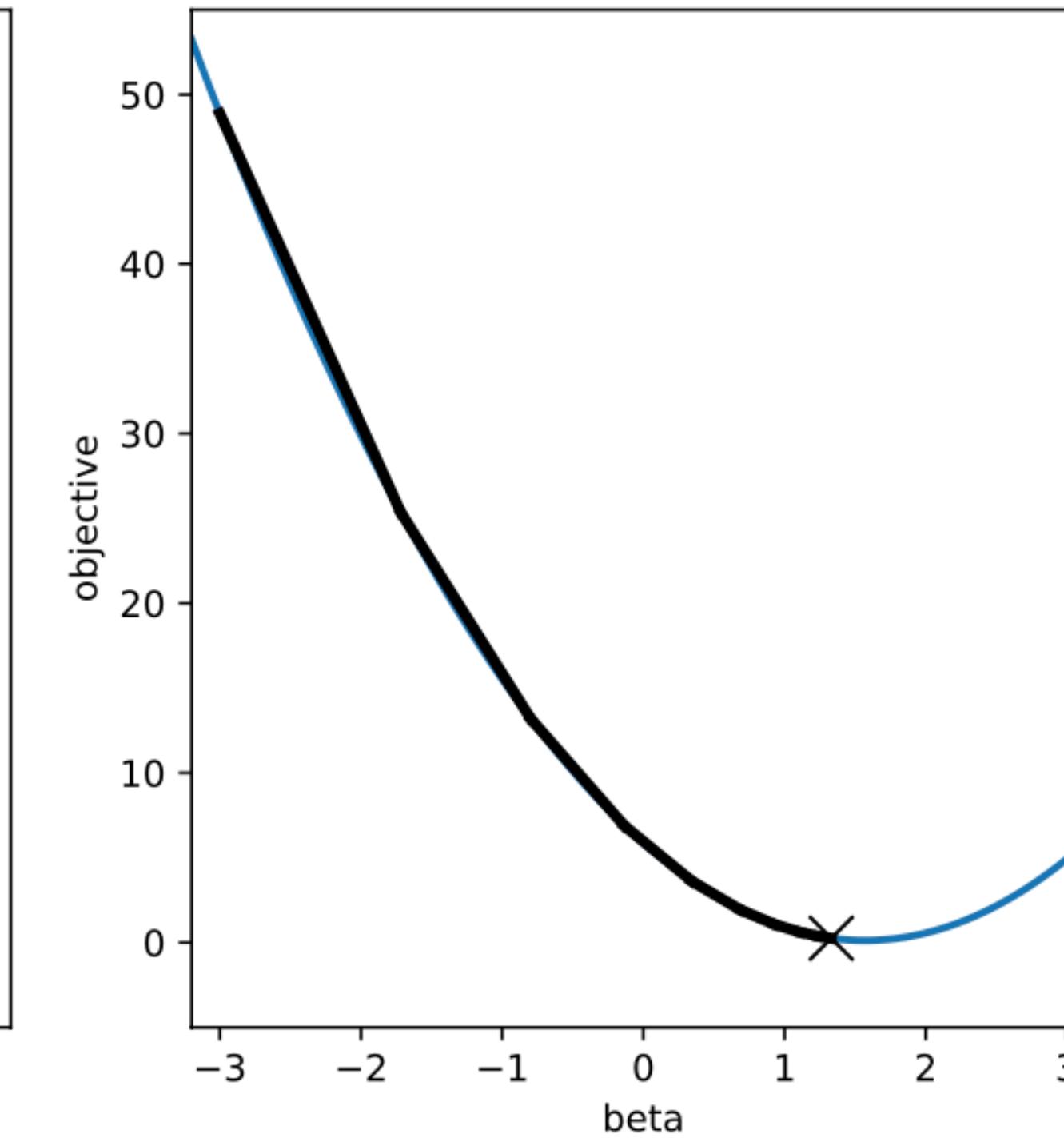
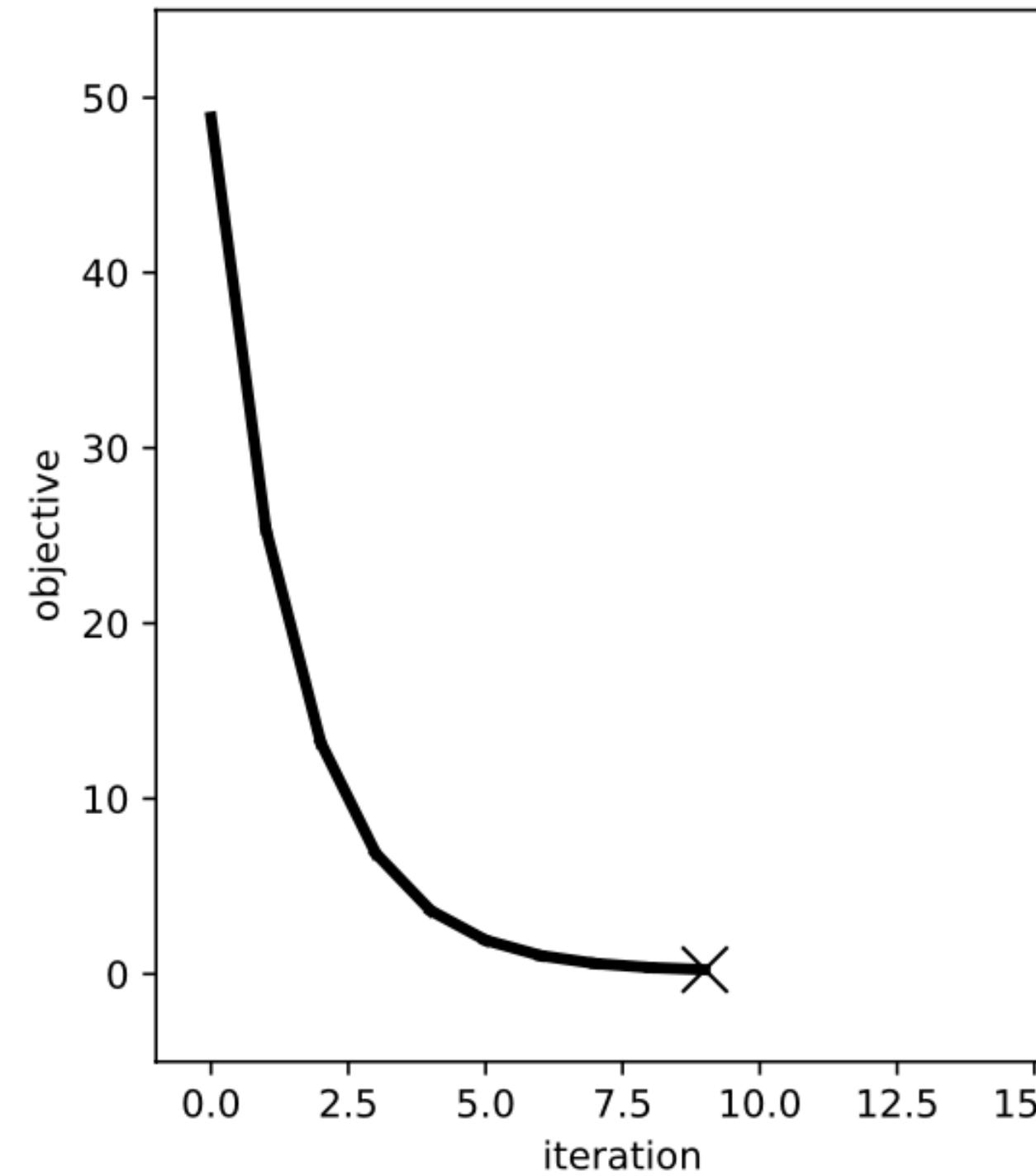
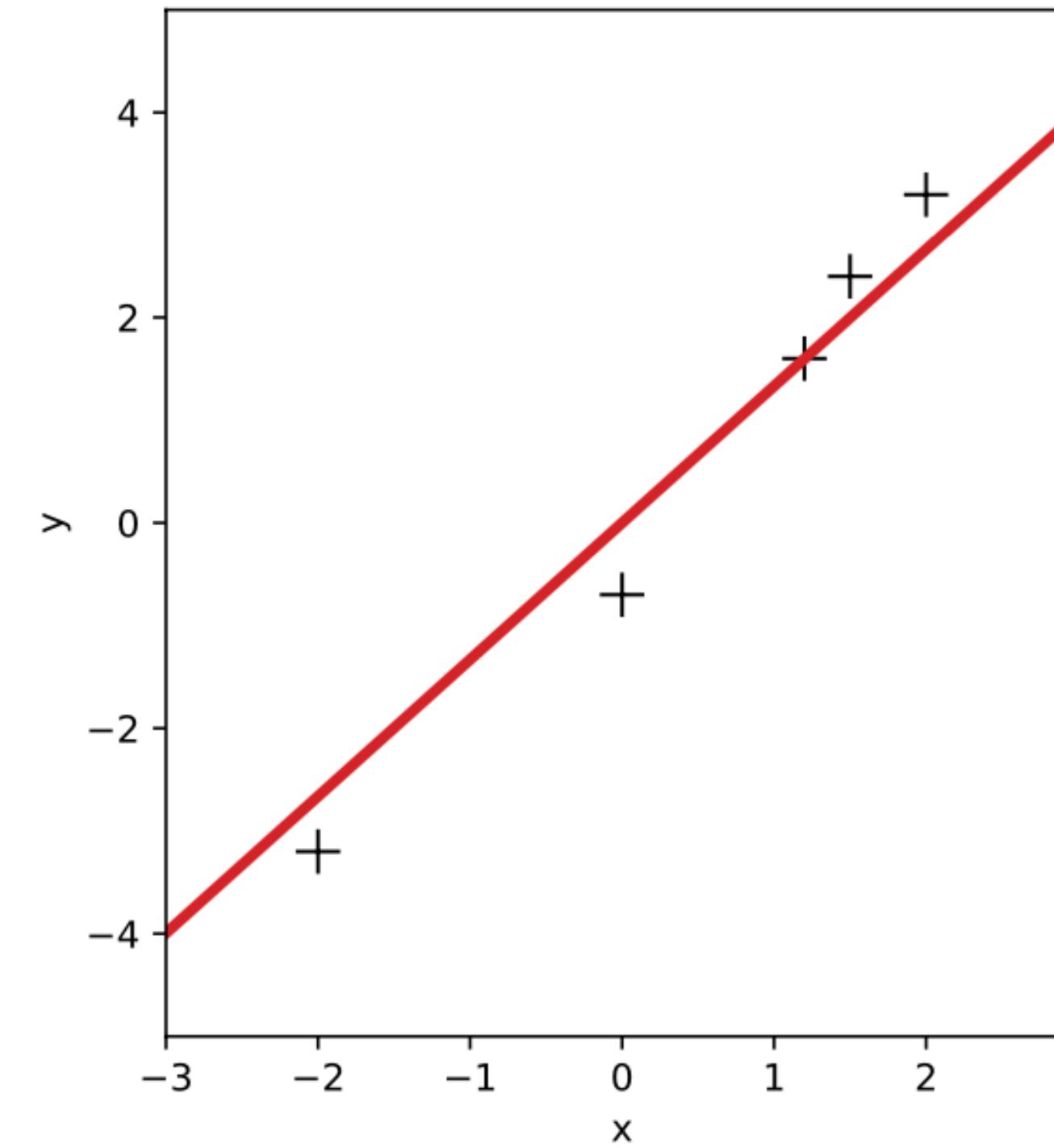
Another example: <https://playground.tensorflow.org/>

# Recap: Optimisation for linear regression



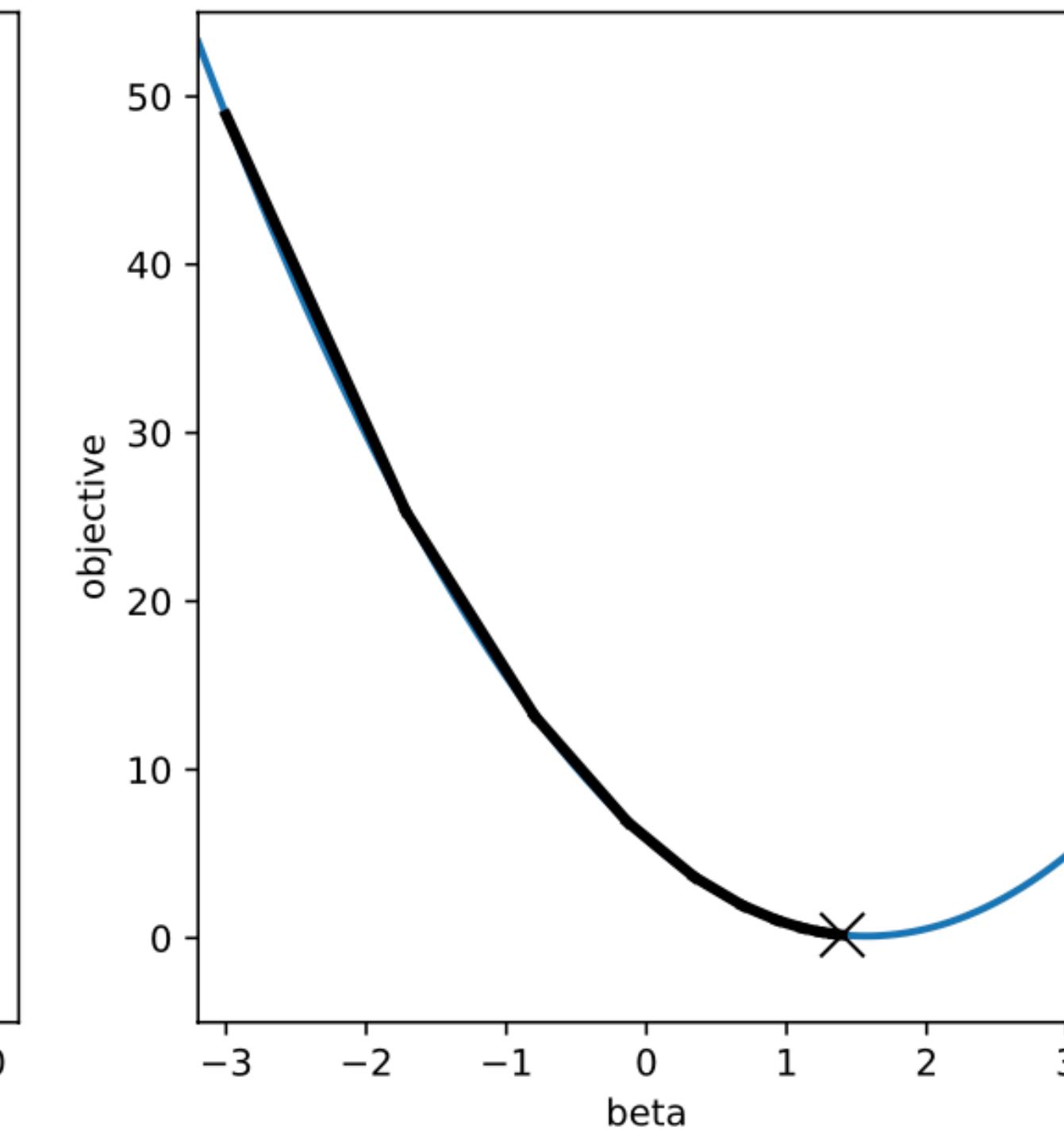
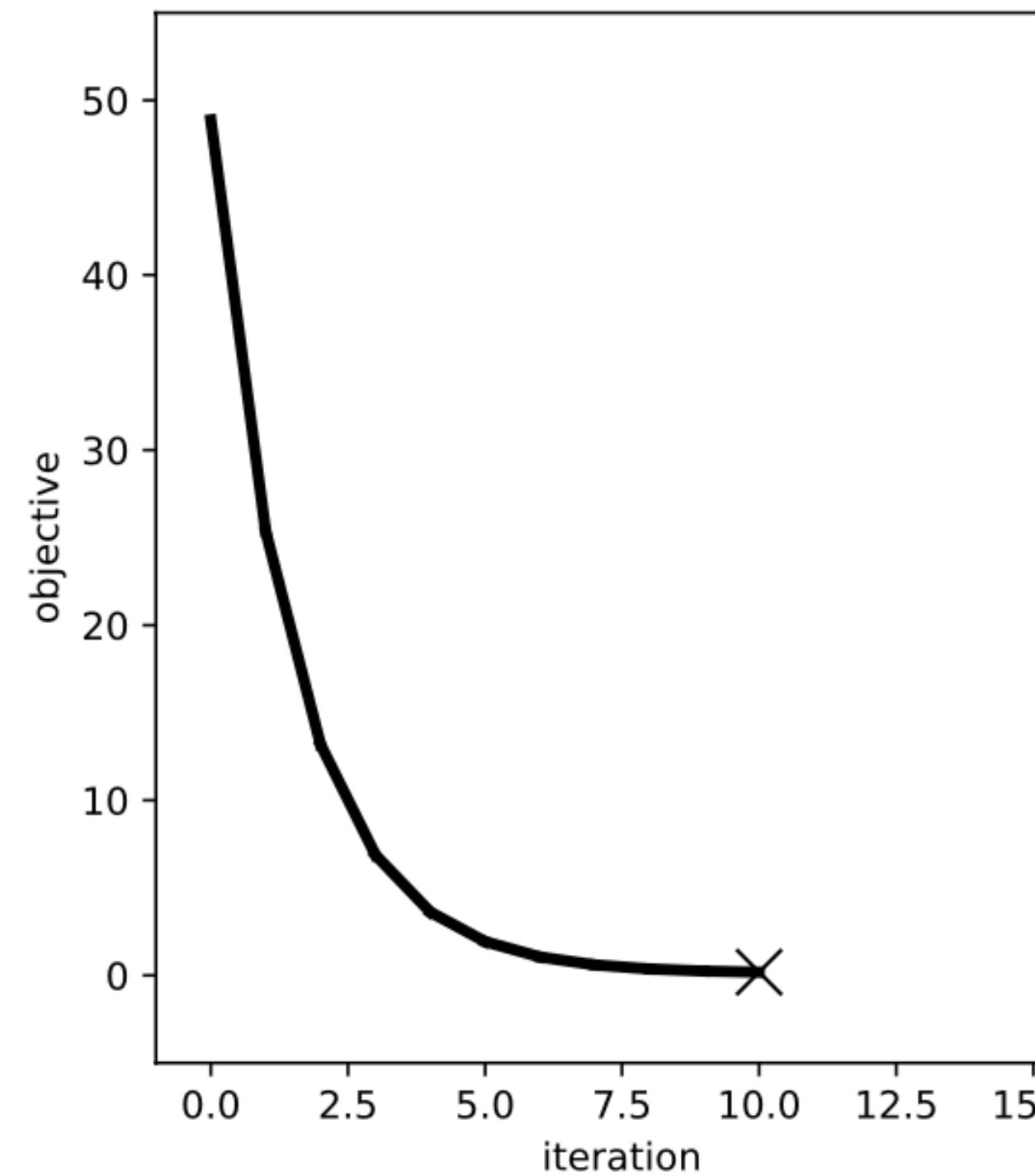
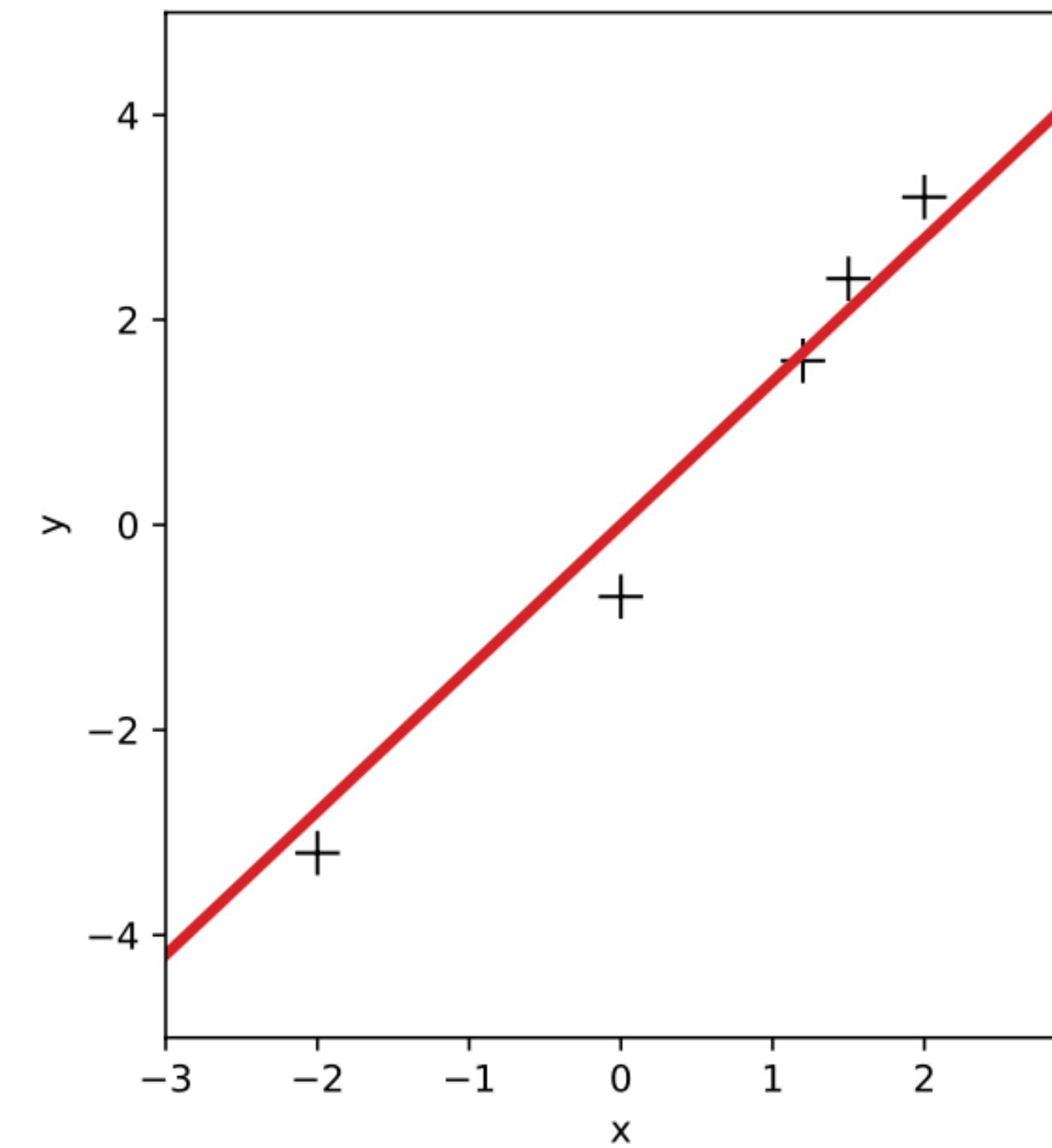
Another example: <https://playground.tensorflow.org/>

# Recap: Optimisation for linear regression



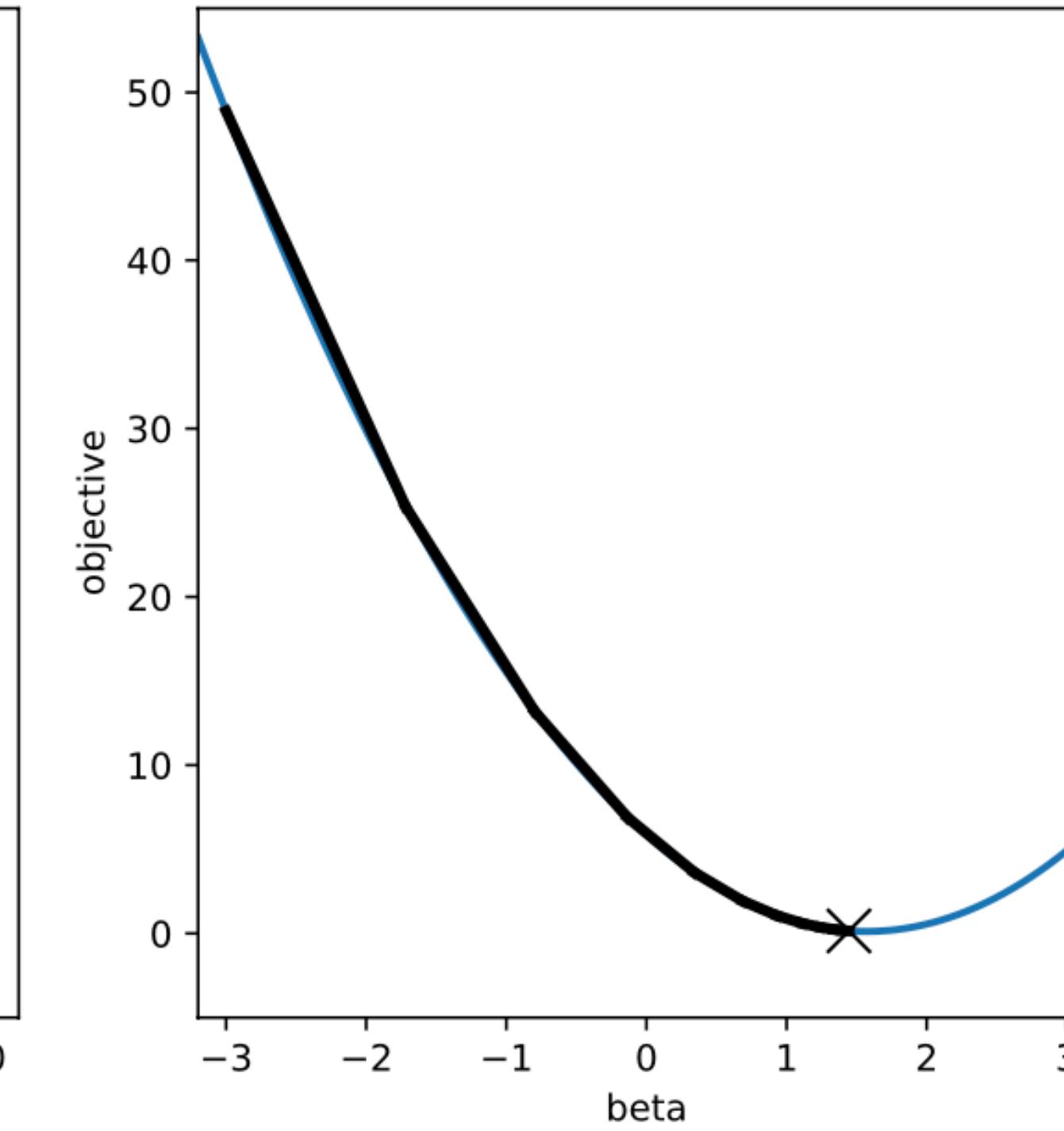
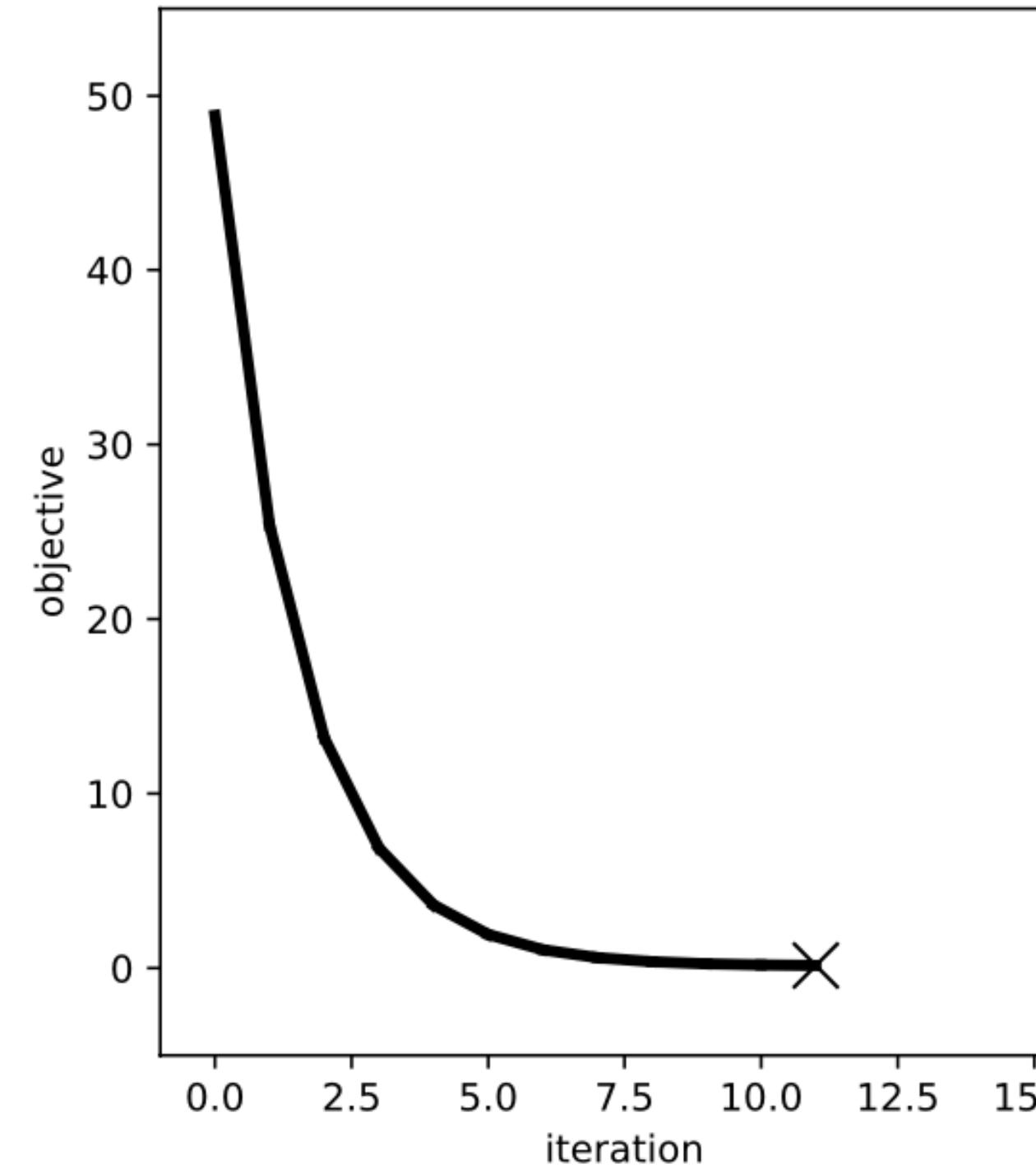
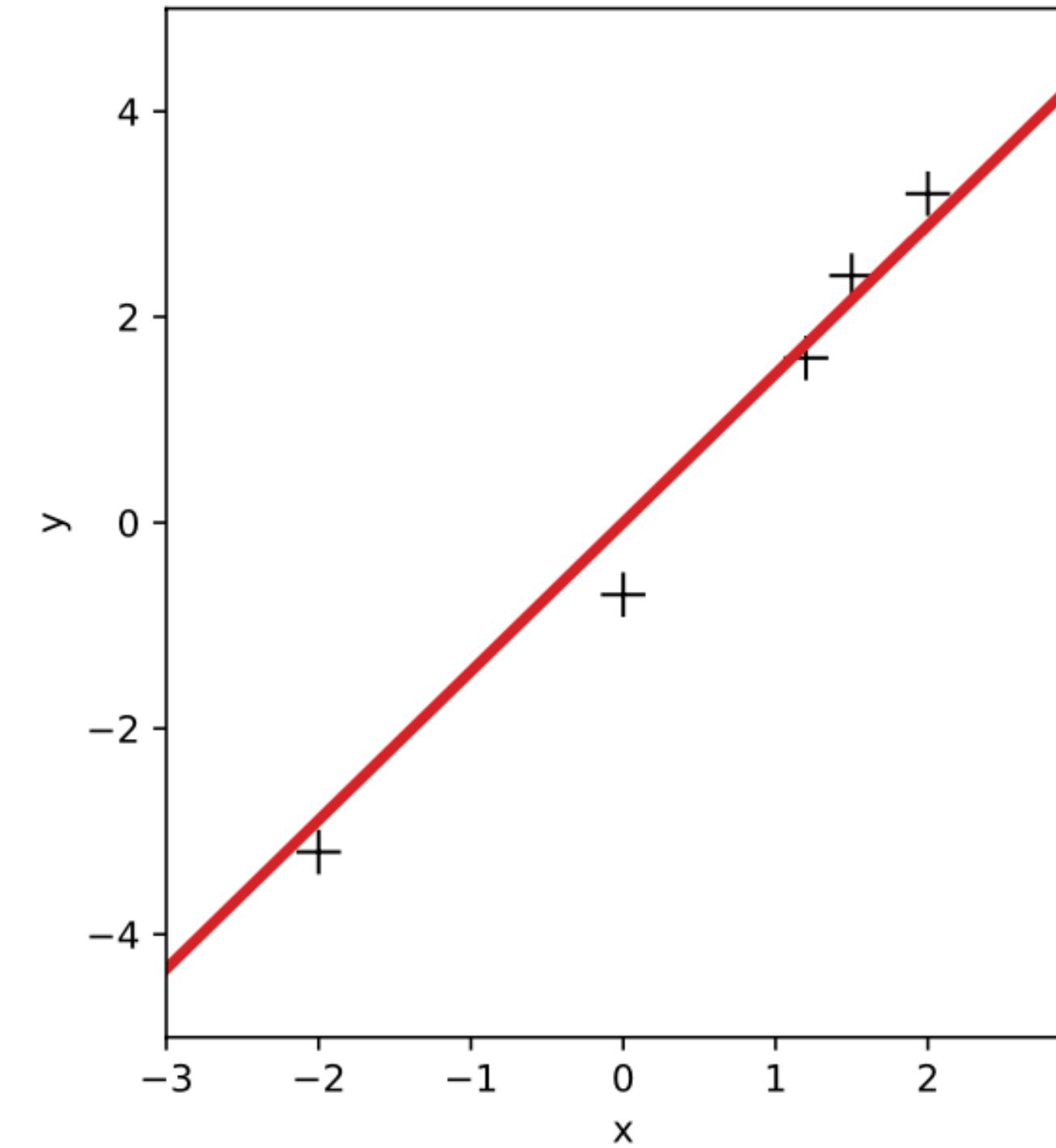
Another example: <https://playground.tensorflow.org/>

# Recap: Optimisation for linear regression



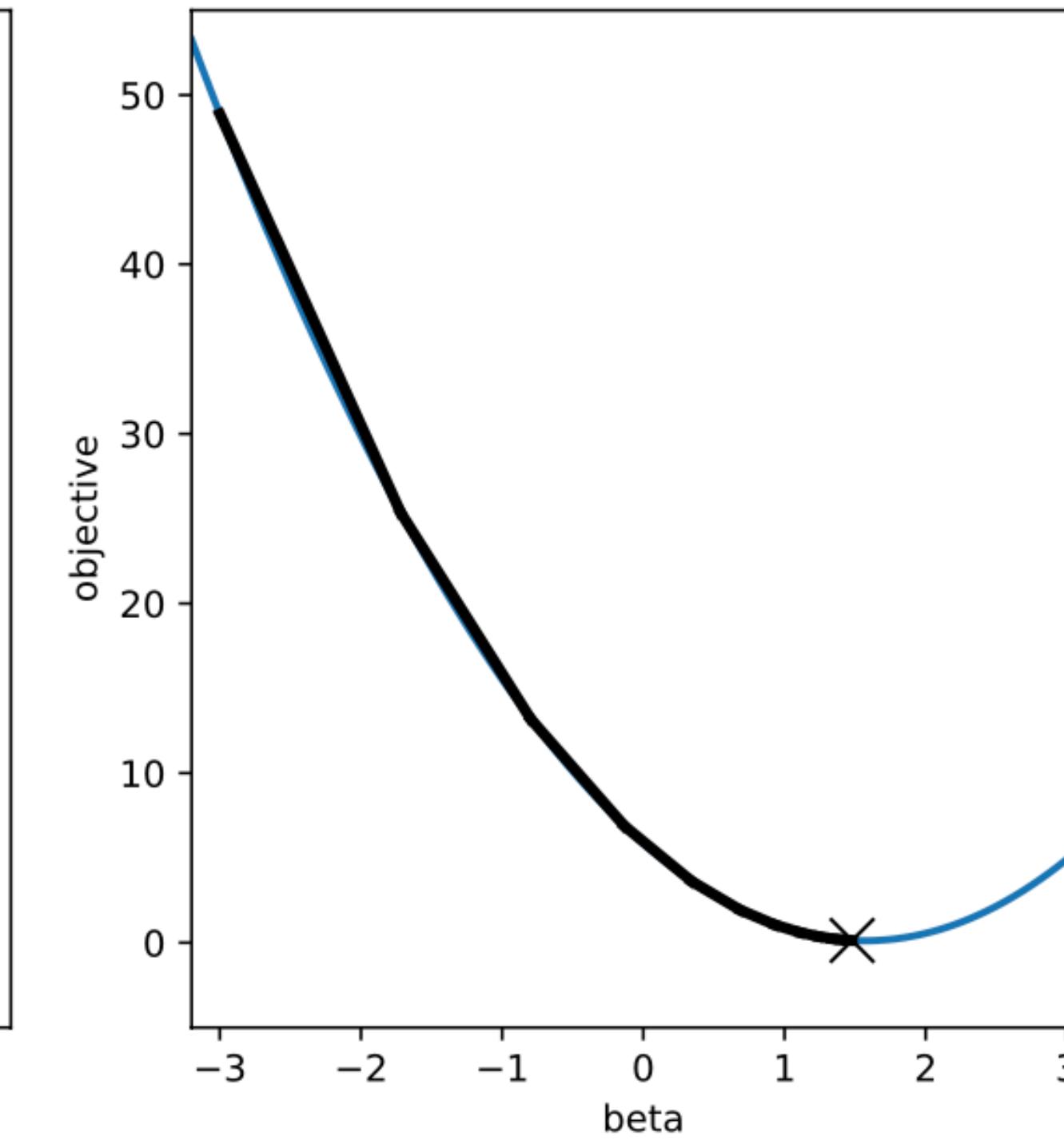
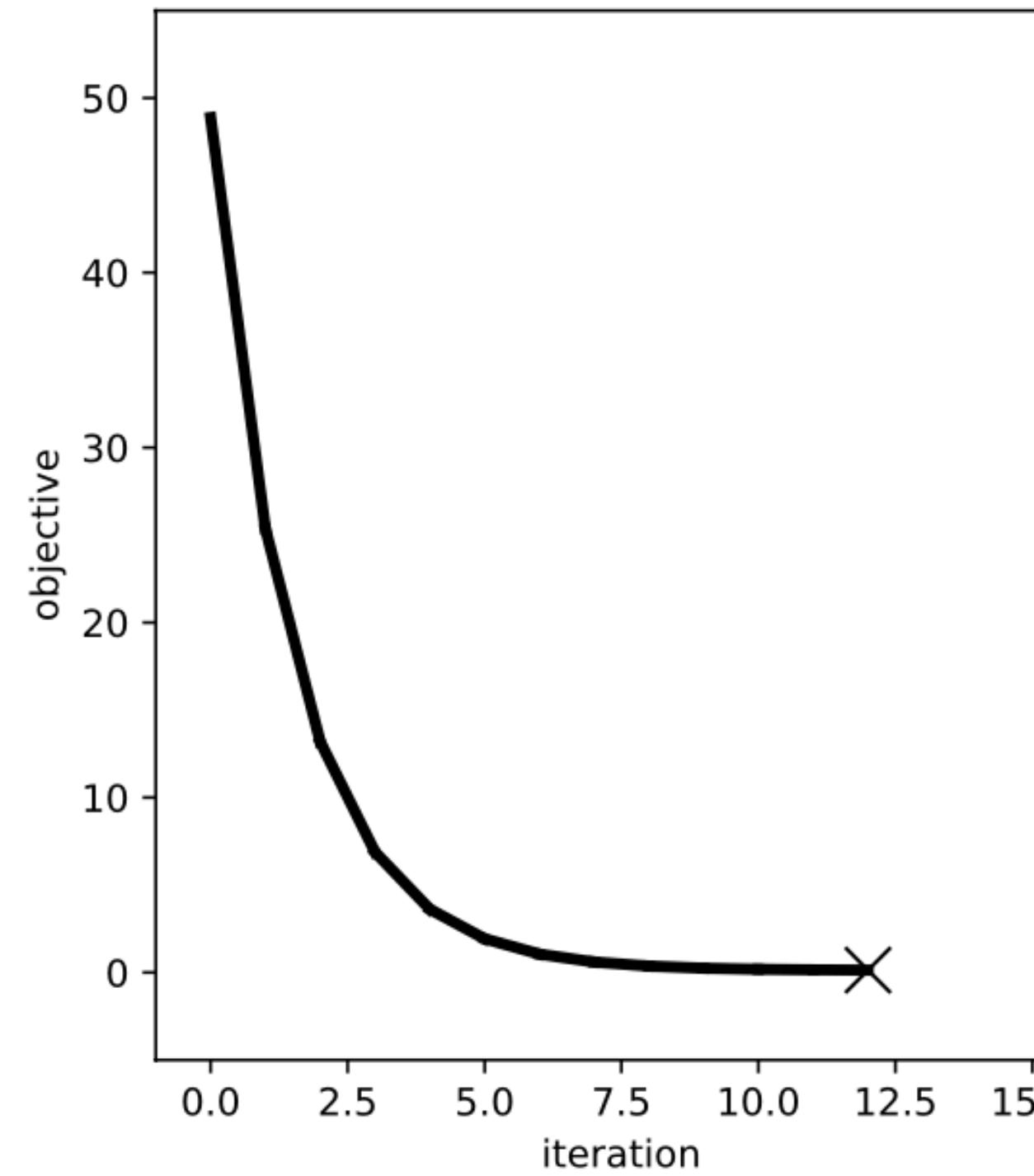
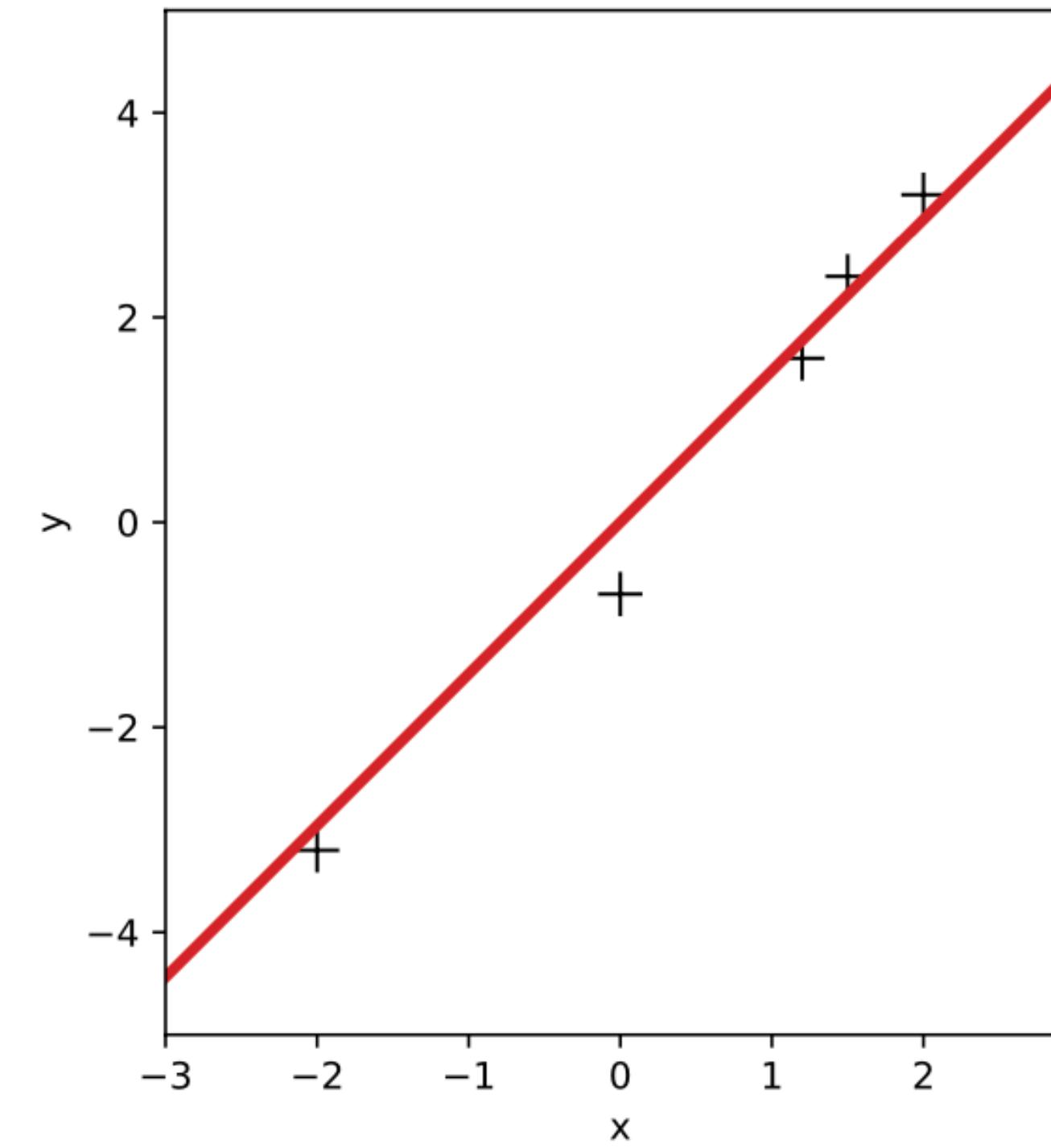
Another example: <https://playground.tensorflow.org/>

# Recap: Optimisation for linear regression



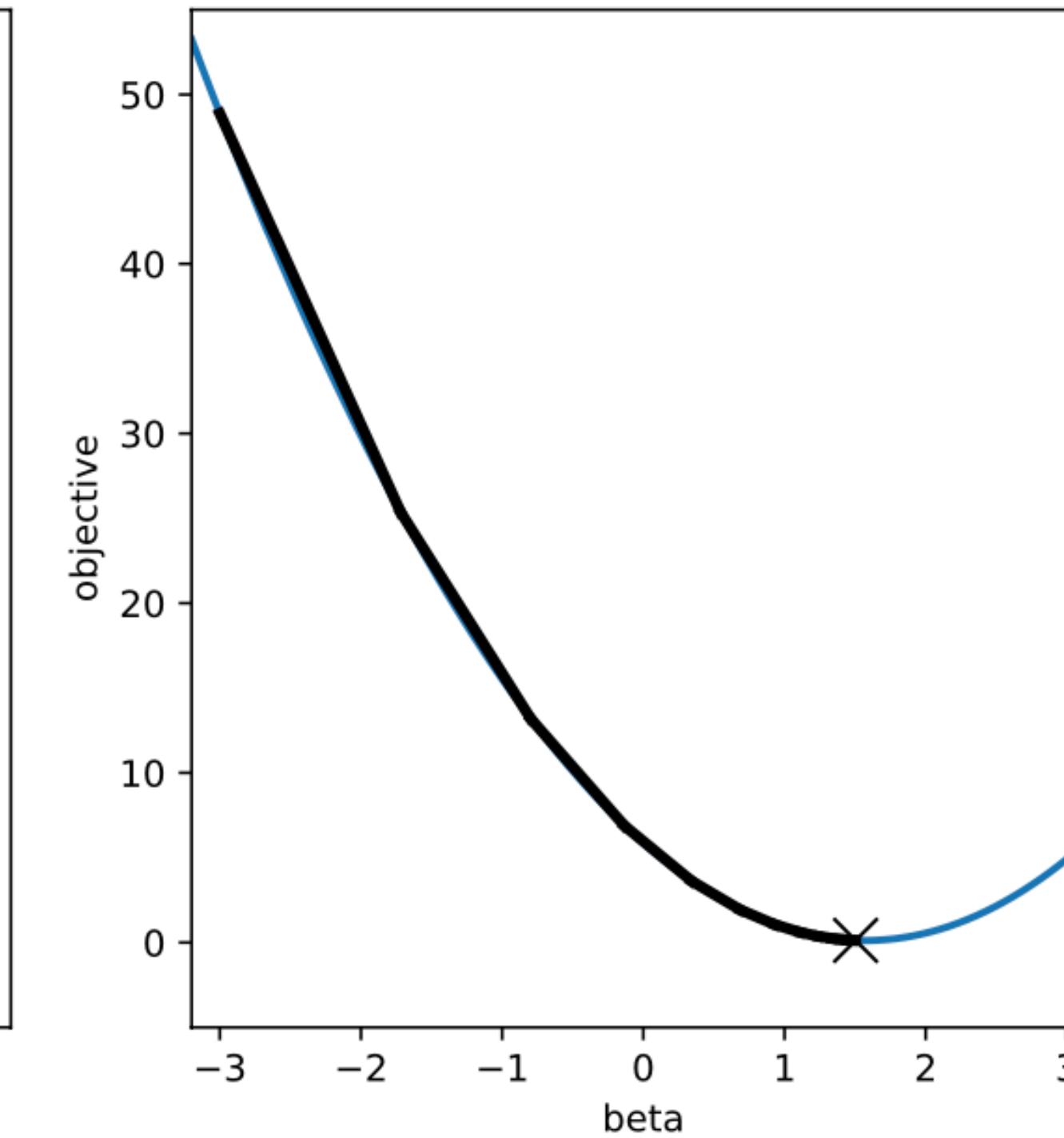
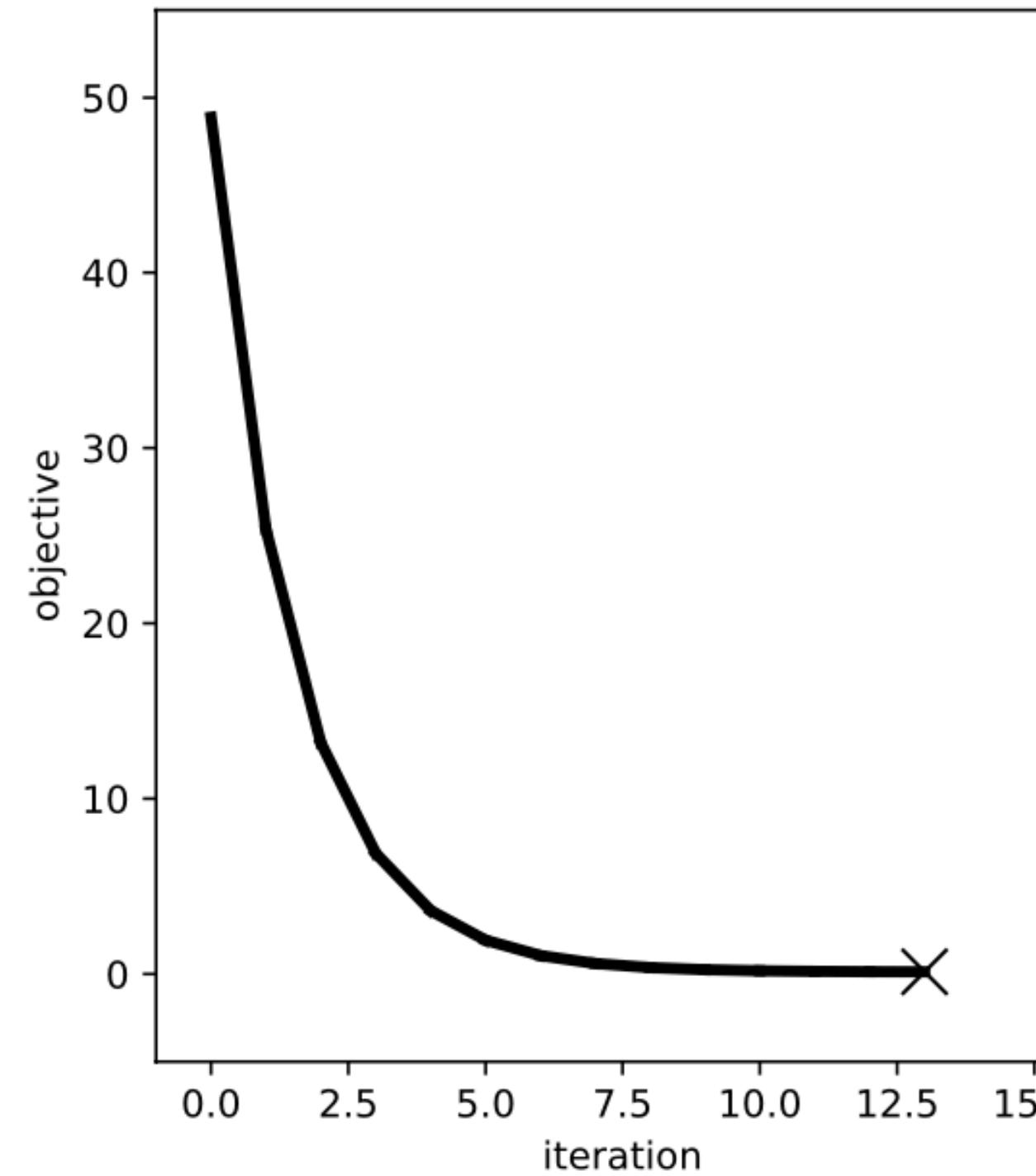
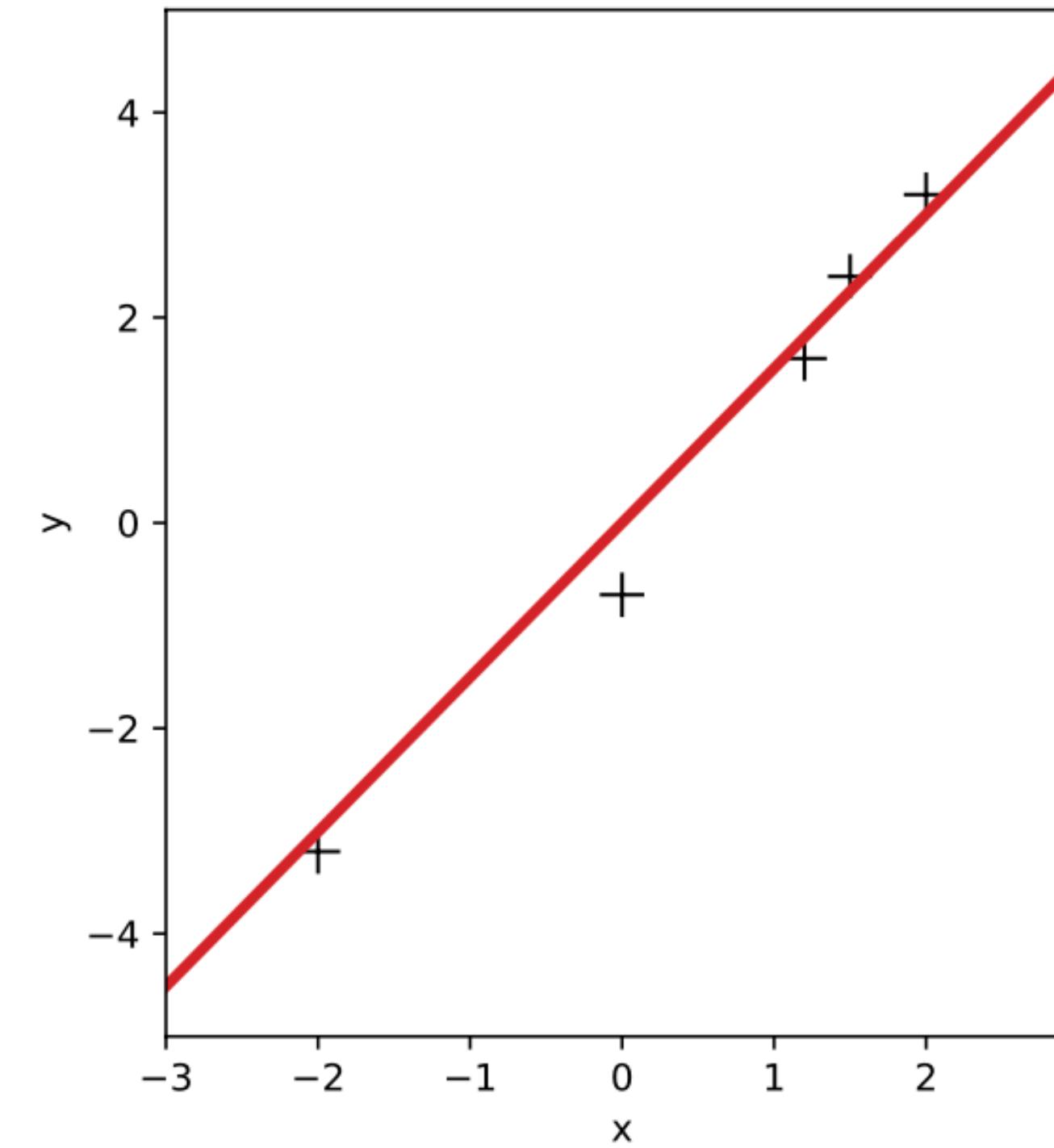
Another example: <https://playground.tensorflow.org/>

# Recap: Optimisation for linear regression



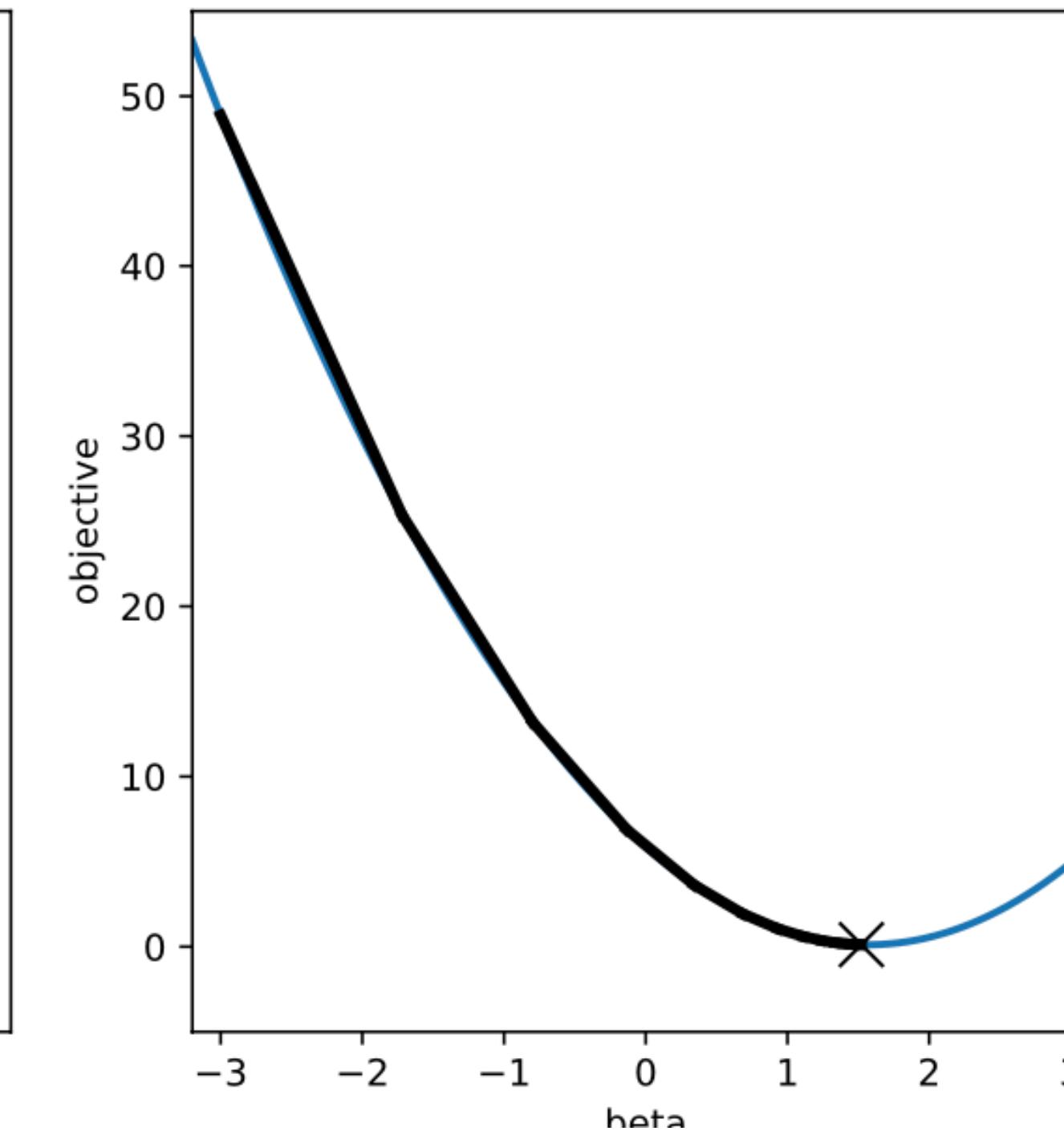
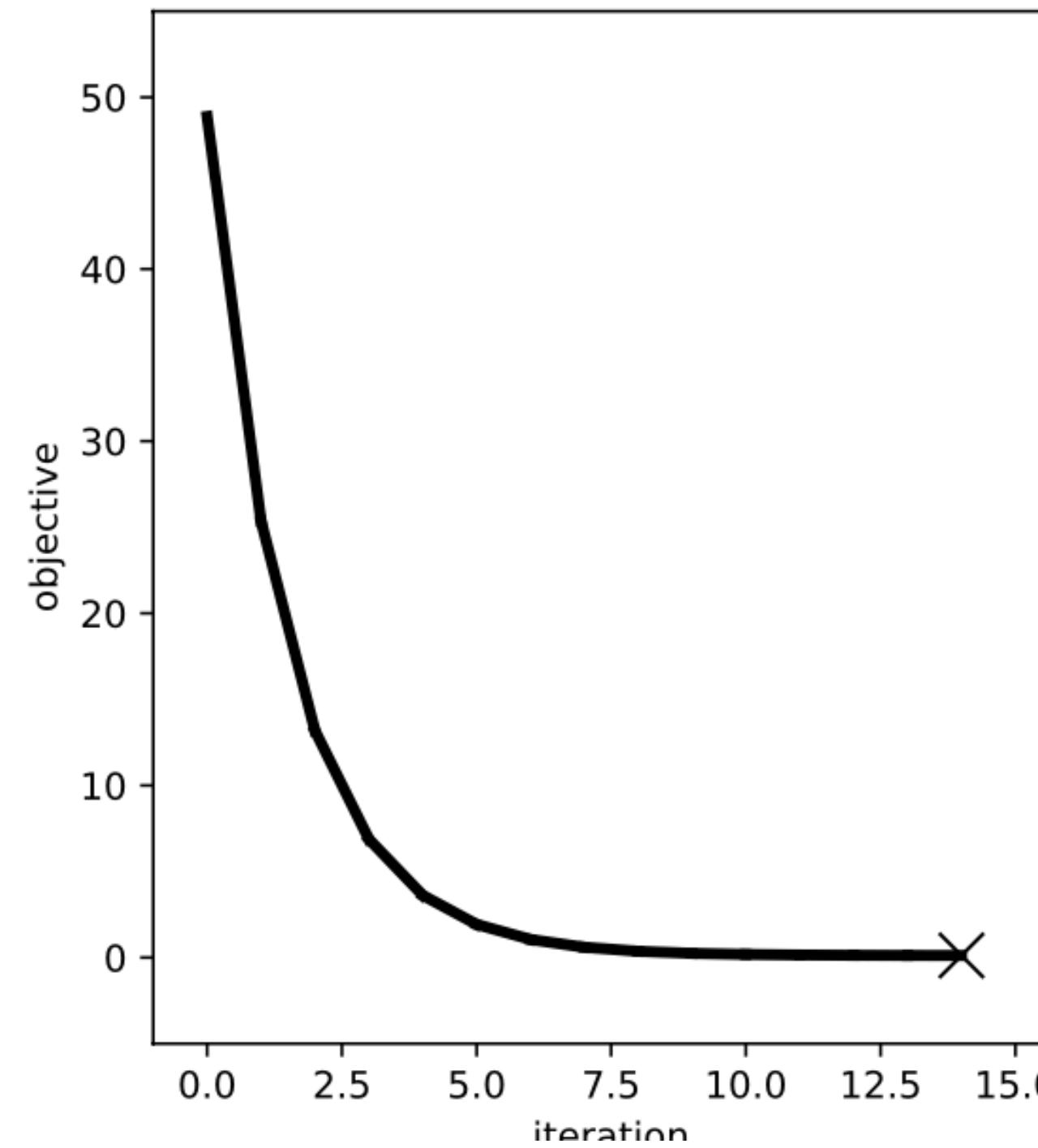
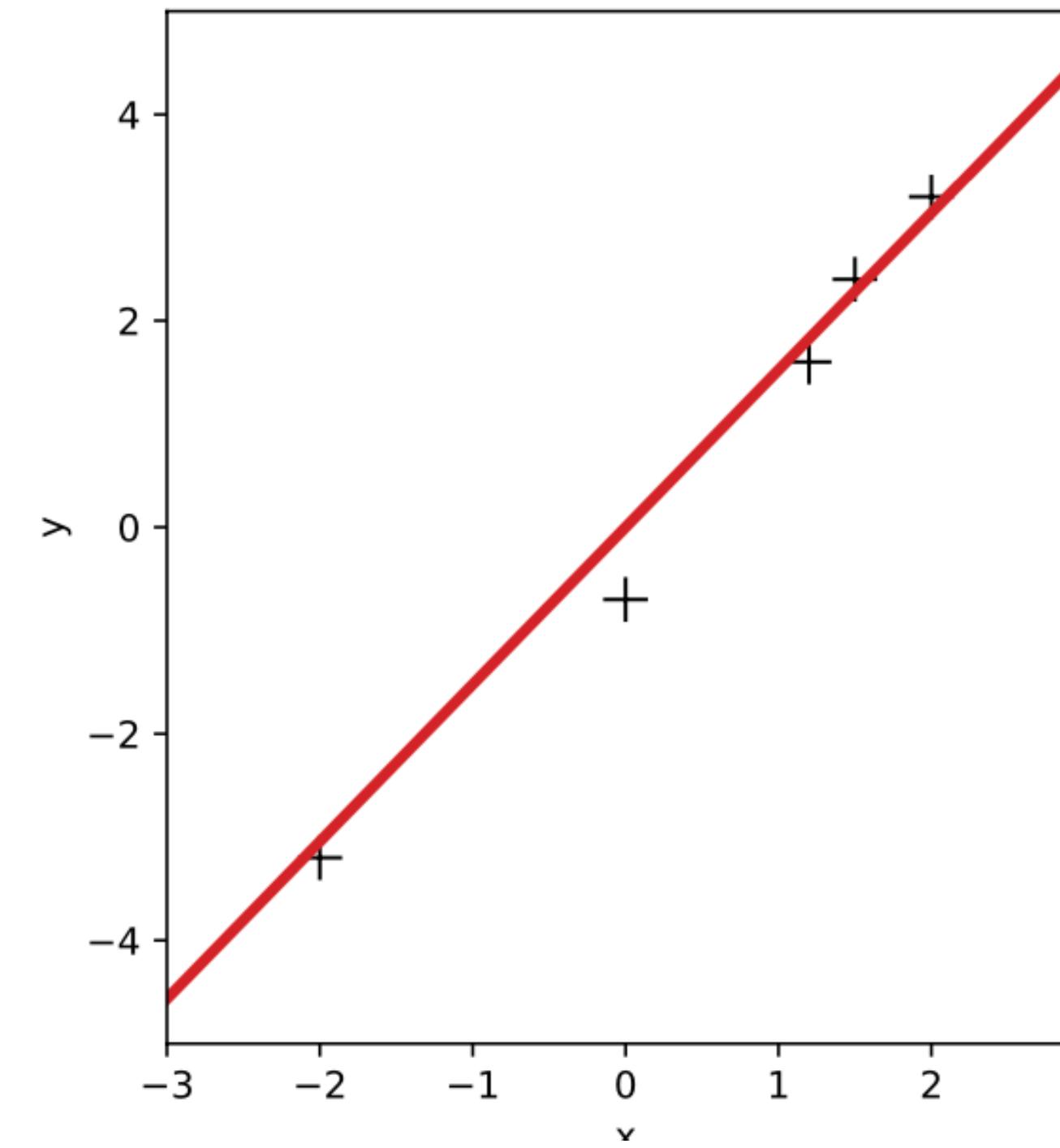
Another example: <https://playground.tensorflow.org/>

# Recap: Optimisation for linear regression



Another example: <https://playground.tensorflow.org/>

# Recap: Optimisation for linear regression



Another example: <https://playground.tensorflow.org/>

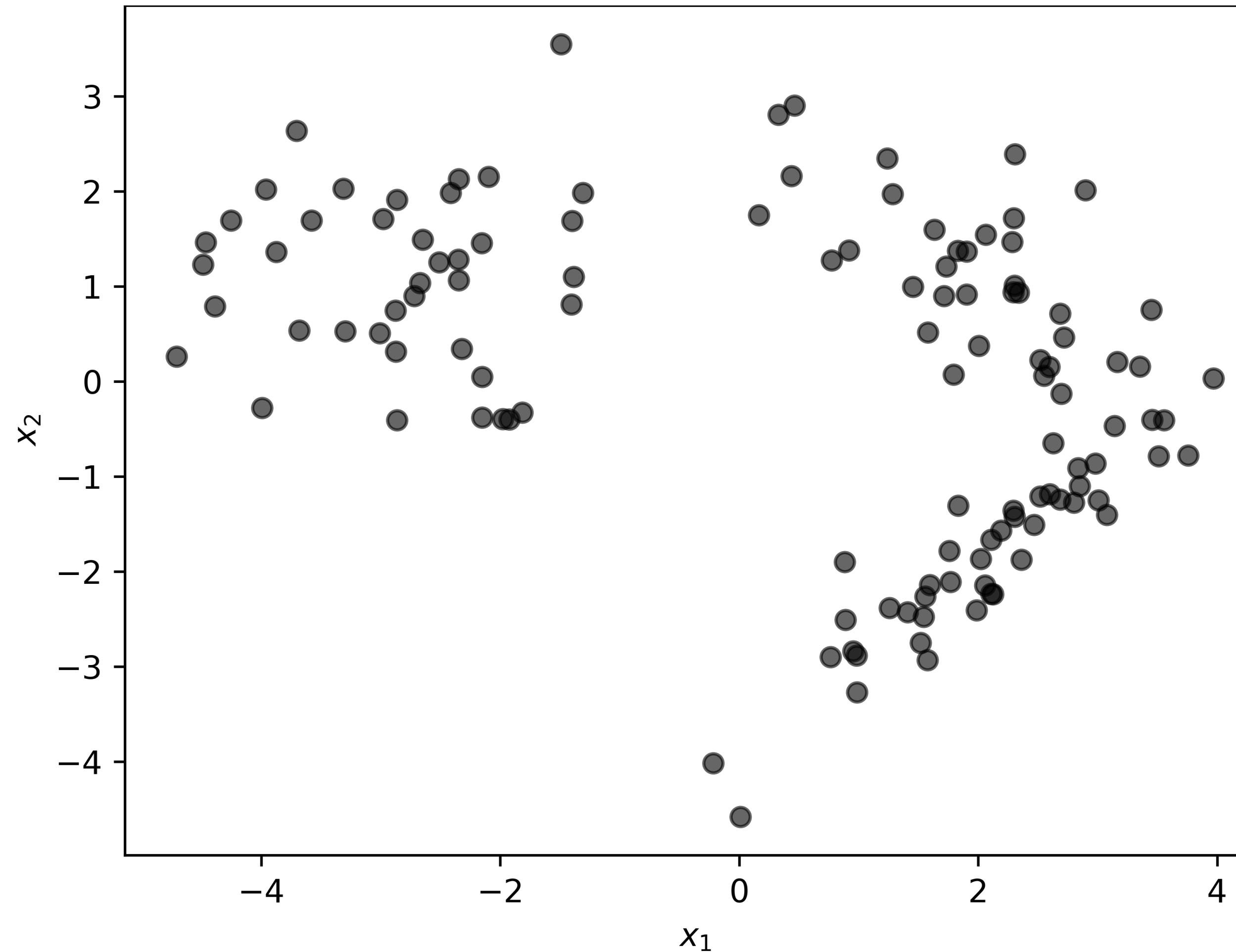
# What is density estimation?

# What is density estimation?

**Goal:** estimate the distribution your data come from, aka estimate  $p(x)$

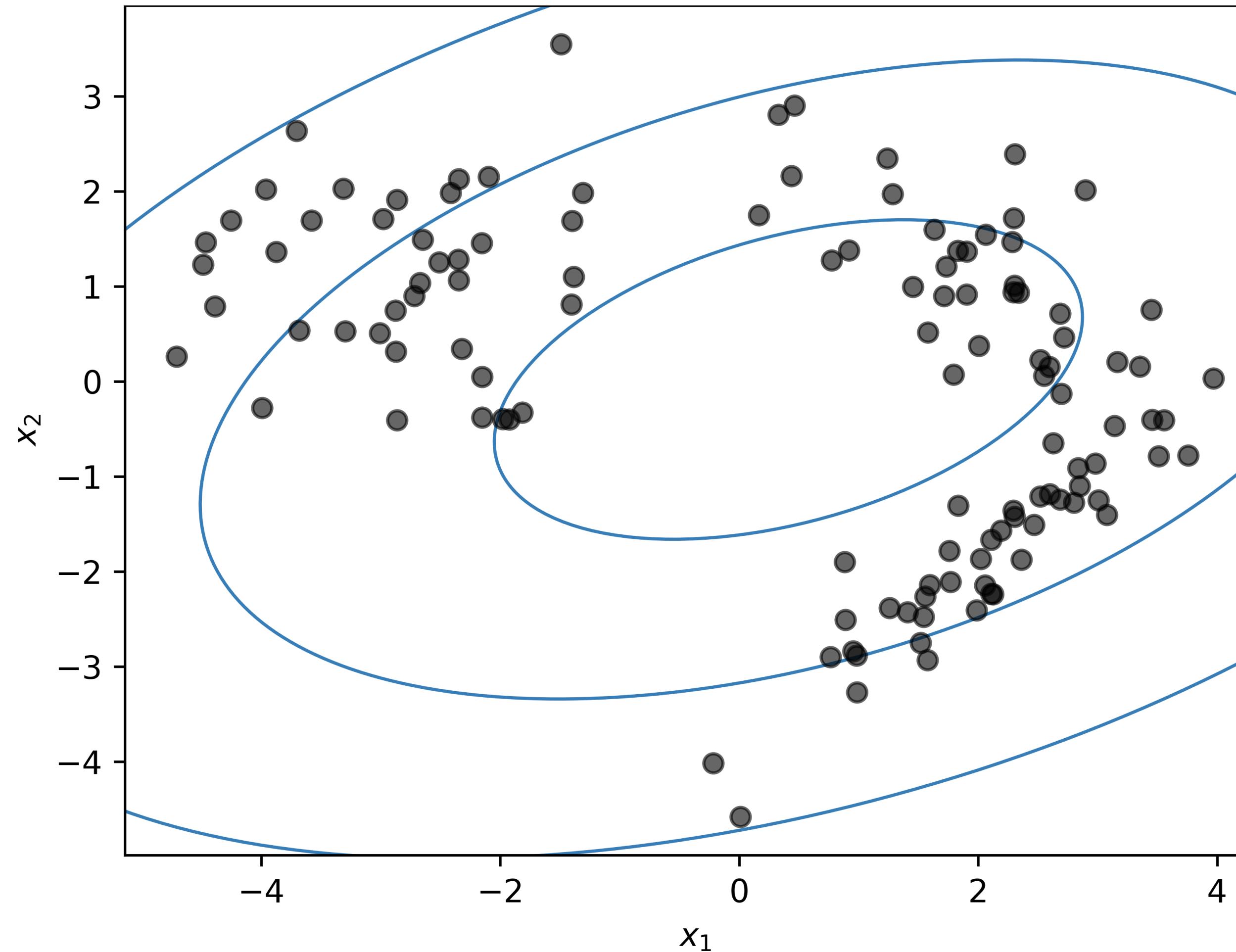
# What is density estimation?

**Goal:** estimate the distribution your data come from, aka estimate  $p(x)$



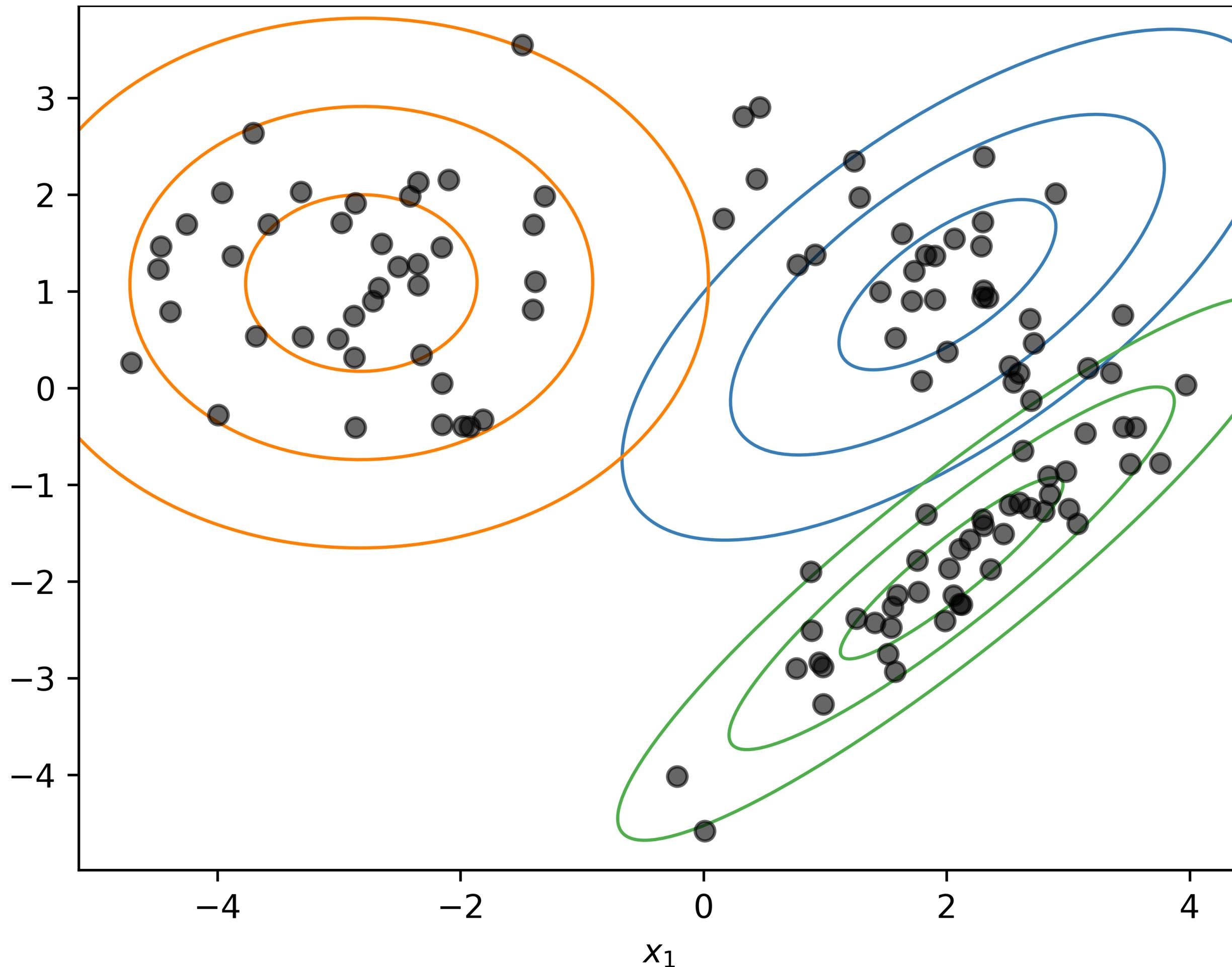
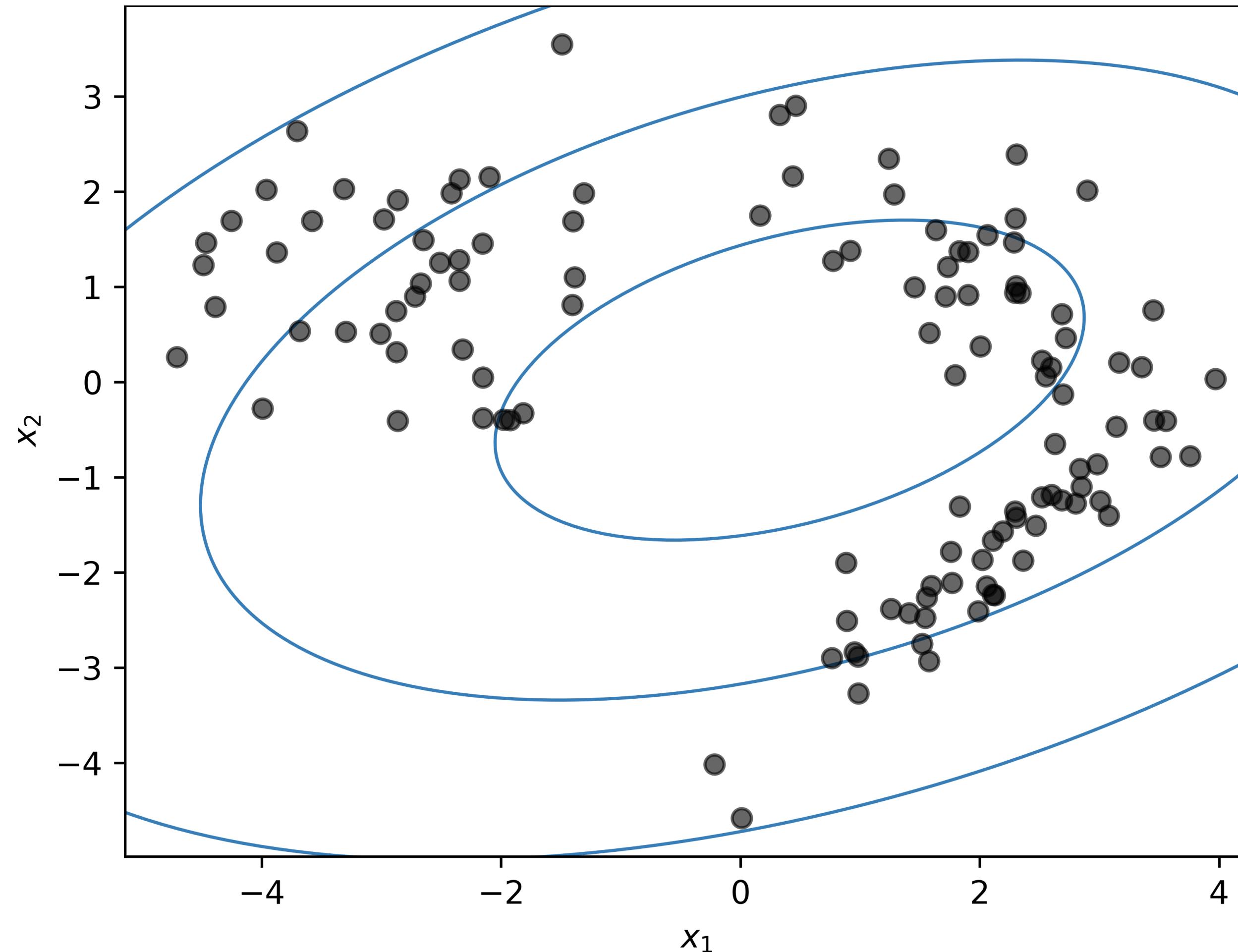
# What is density estimation?

**Goal:** estimate the distribution your data come from, aka estimate  $p(x)$



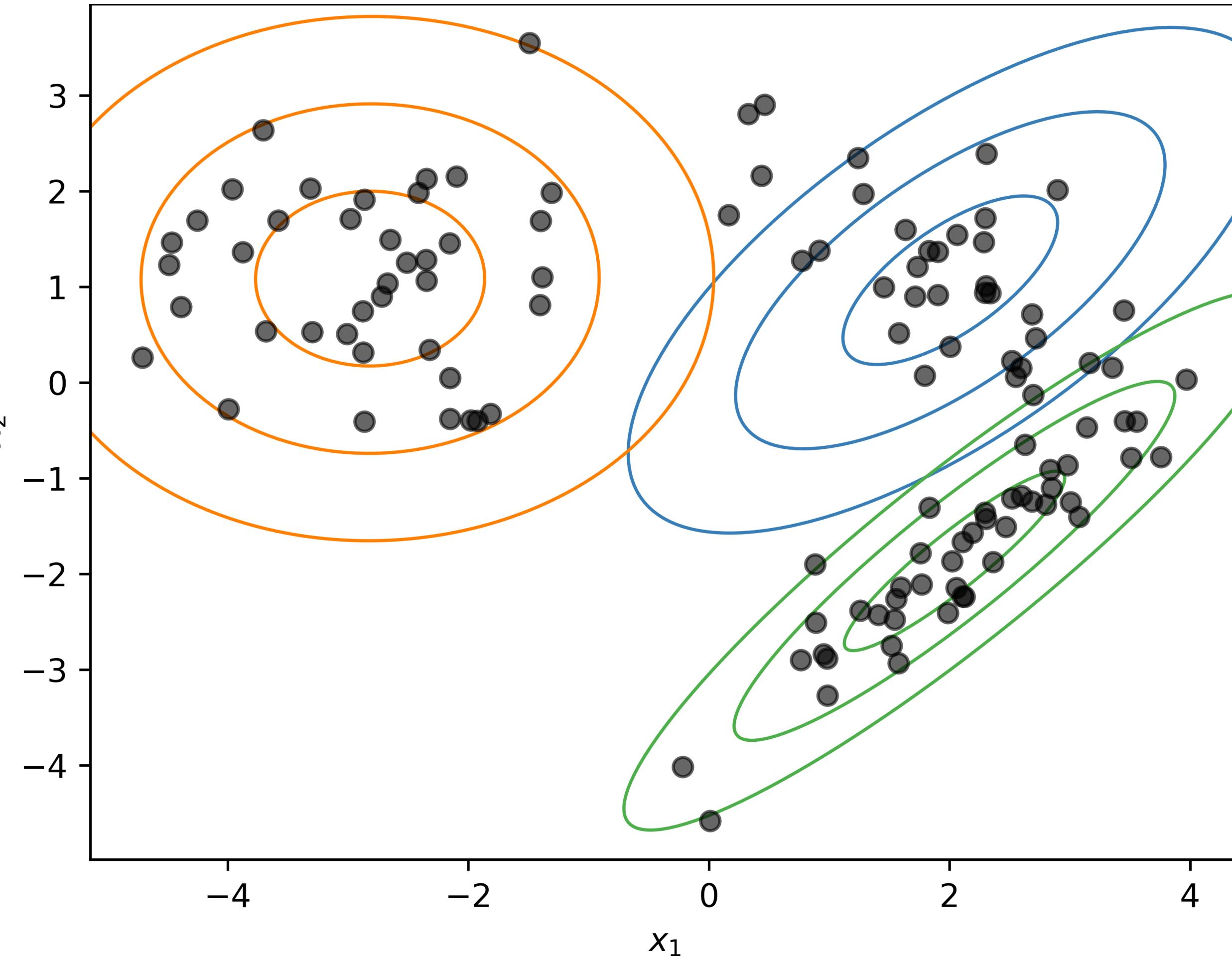
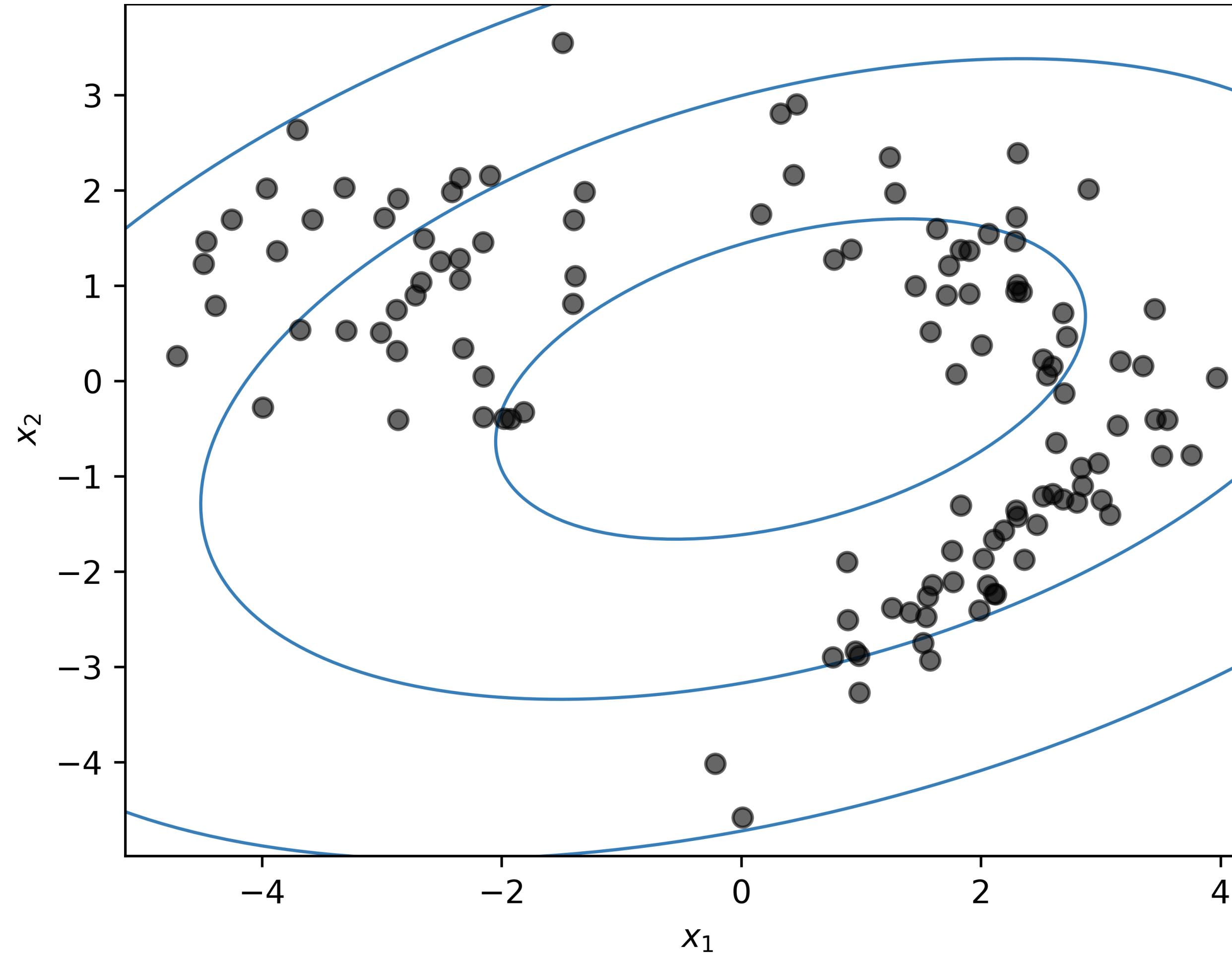
# What is density estimation?

Goal: estimate the distribution your data come from, aka estimate  $p(x)$



# What is density estimation?

Goal: estimate the distribution your data come from, aka estimate  $p(x)$



Applications: (1) exploratory analyses, (2) evaluation, (3) synthetic data generation

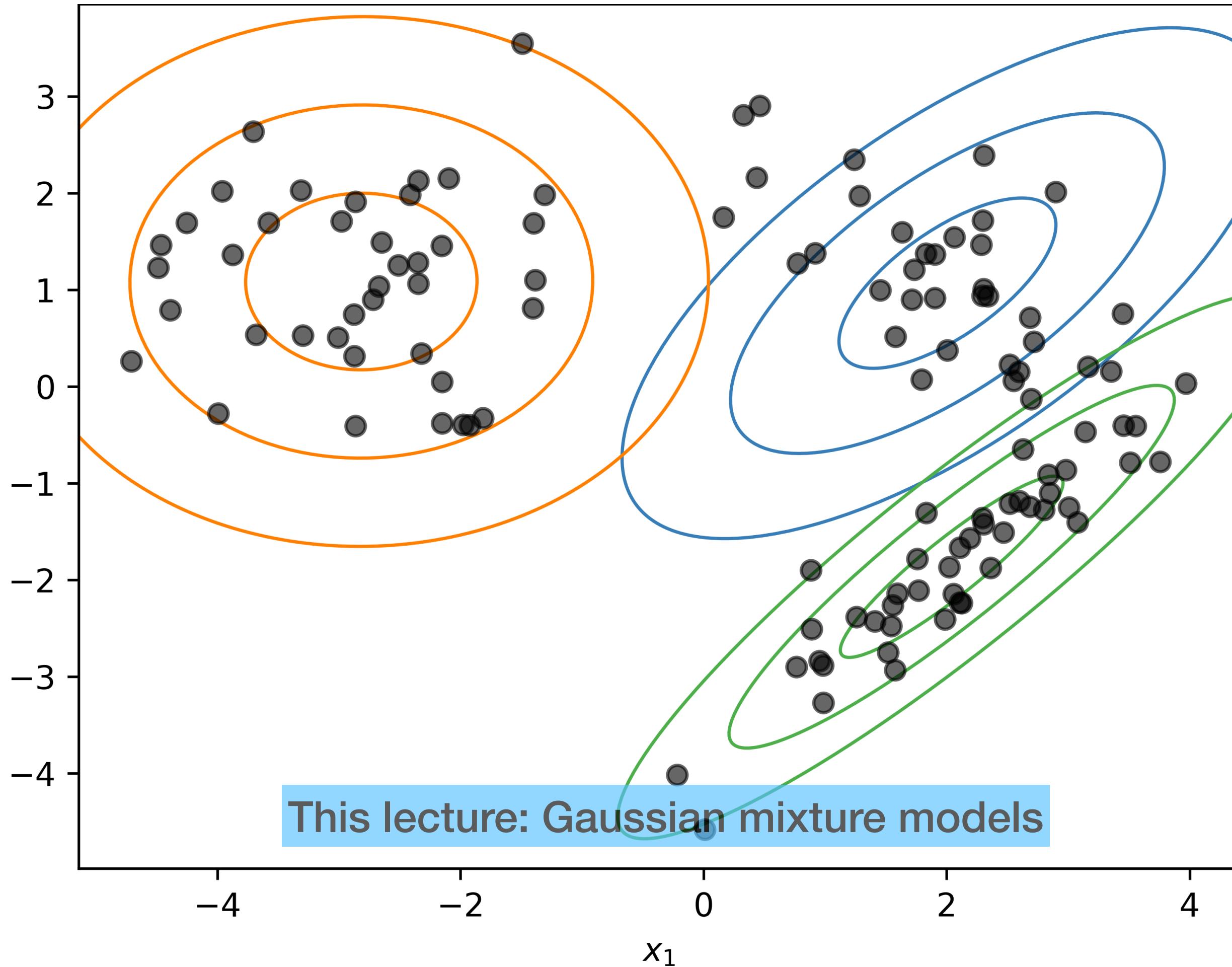
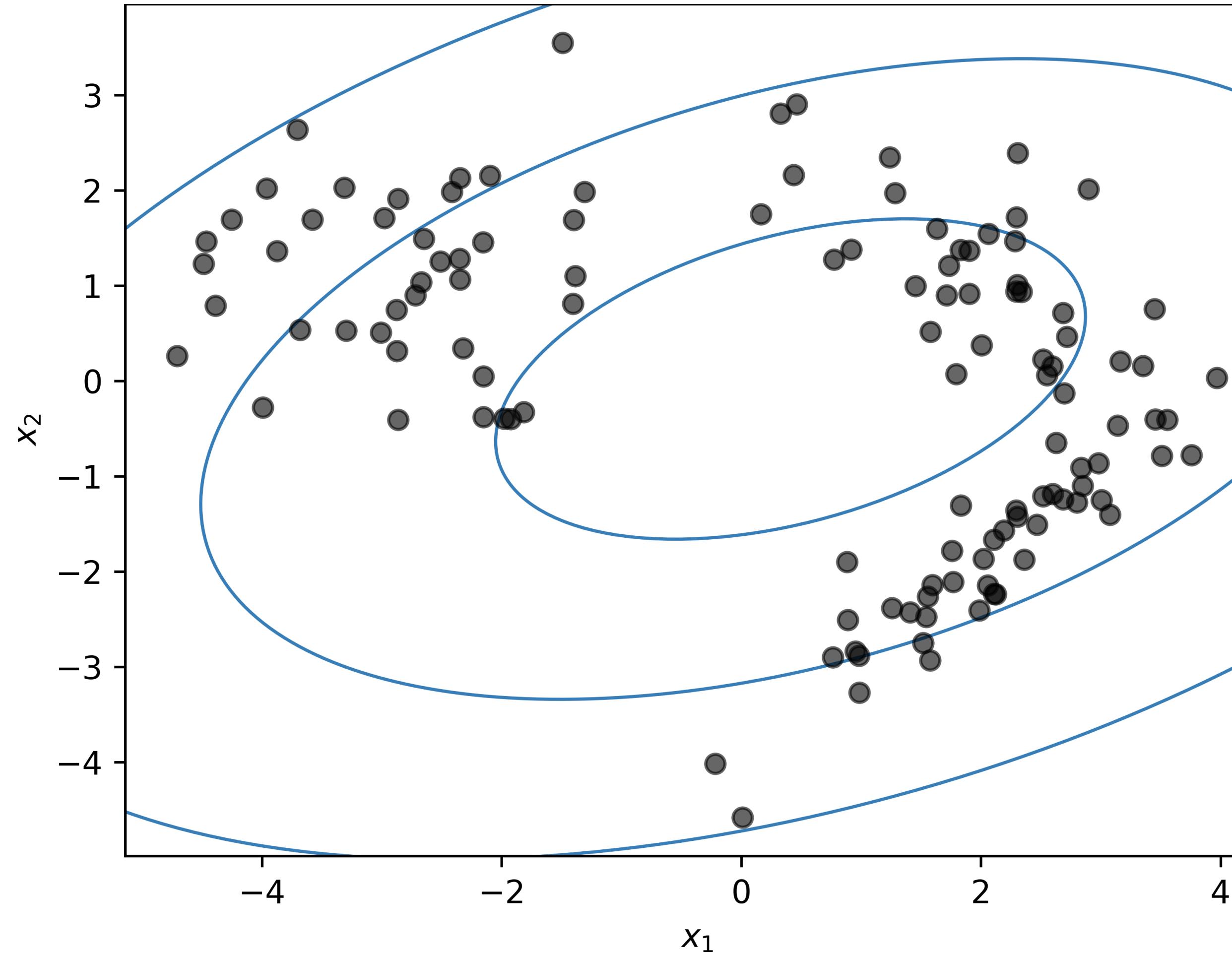
Are the data bimodal or skewed?

Do datapoints come from the same density?

Generate new samples that have the same characteristics

# What is density estimation?

Goal: estimate the distribution your data come from, aka estimate  $p(x)$



Applications: (1) exploratory analyses, (2) evaluation, (3) synthetic data generation

Are the data bimodal or skewed?

Do datapoints come from the same density?

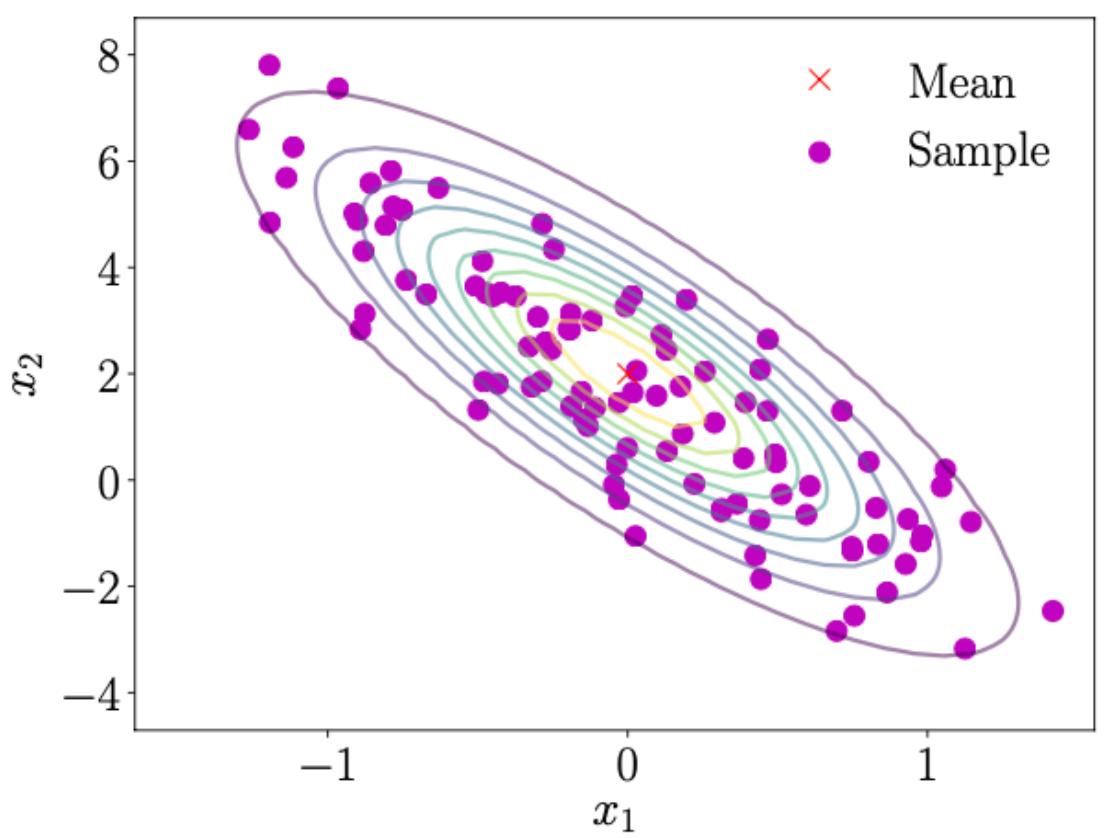
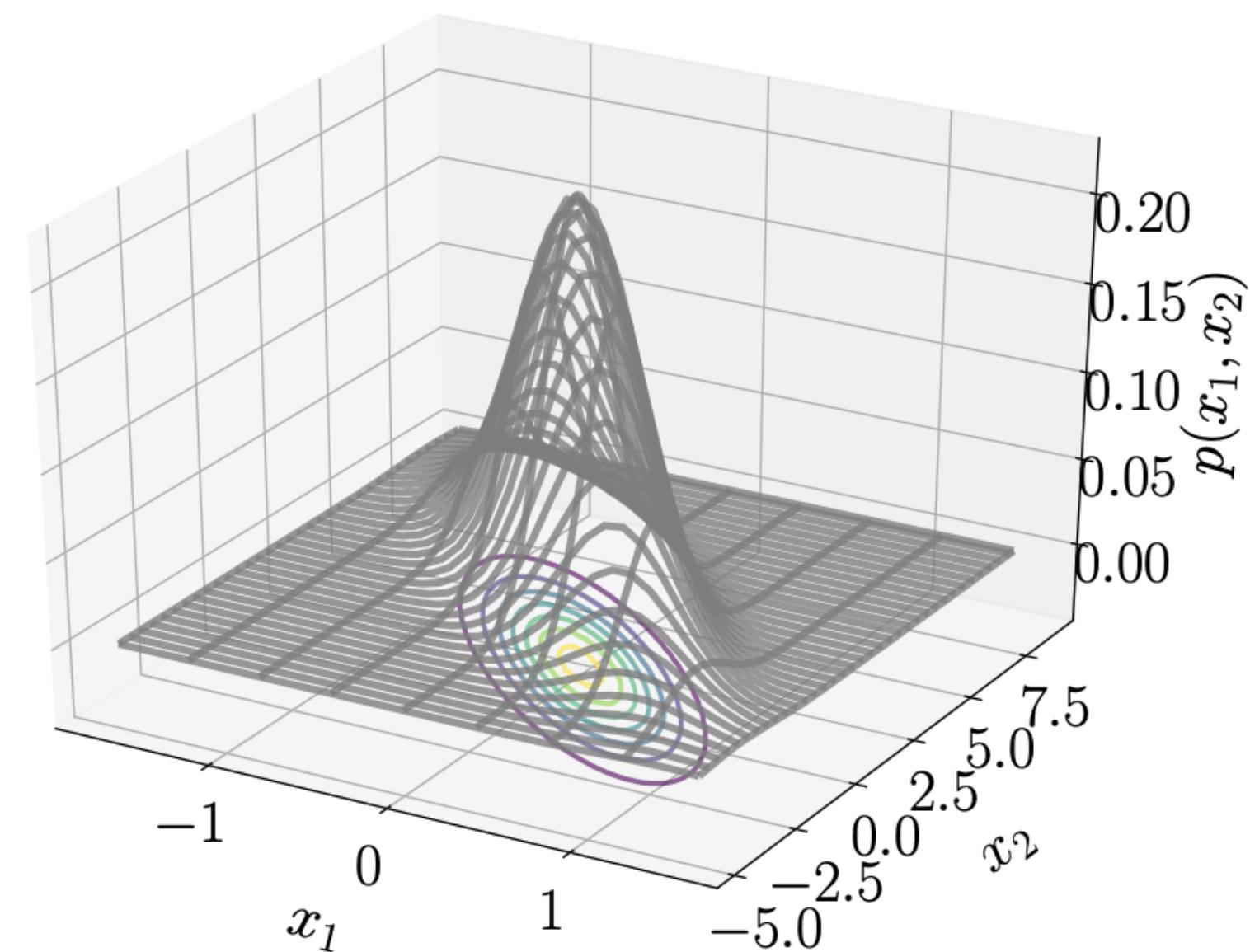
Generate new samples that have the same characteristics

# Recap: multivariate Gaussian distributions (W6)

A D-dimensional multivariate **Gaussian** or **normal** distribution is characterised by a *mean* vector  $\mu$  and a *covariance matrix*  $\Sigma$ , with

$$\text{pdf: } p(\mathbf{x} | \mu, \Sigma) = (2\pi)^{-\frac{D}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right)$$

We often write:  $p(\mathbf{x} | \mu, \Sigma) = \mathcal{N}(\mathbf{x}; \mu, \Sigma)$  or  $\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)$



(b) Multivariate (two-dimensional) Gaussian, viewed from top. The red cross shows the mean and the colored lines show the contour lines of the density.

# Recap: multivariate Gaussian distributions (W6)

A D-dimensional multivariate **Gaussian** or **normal** distribution is characterised by a *mean* vector  $\mu$  and a *covariance matrix*  $\Sigma$ , with

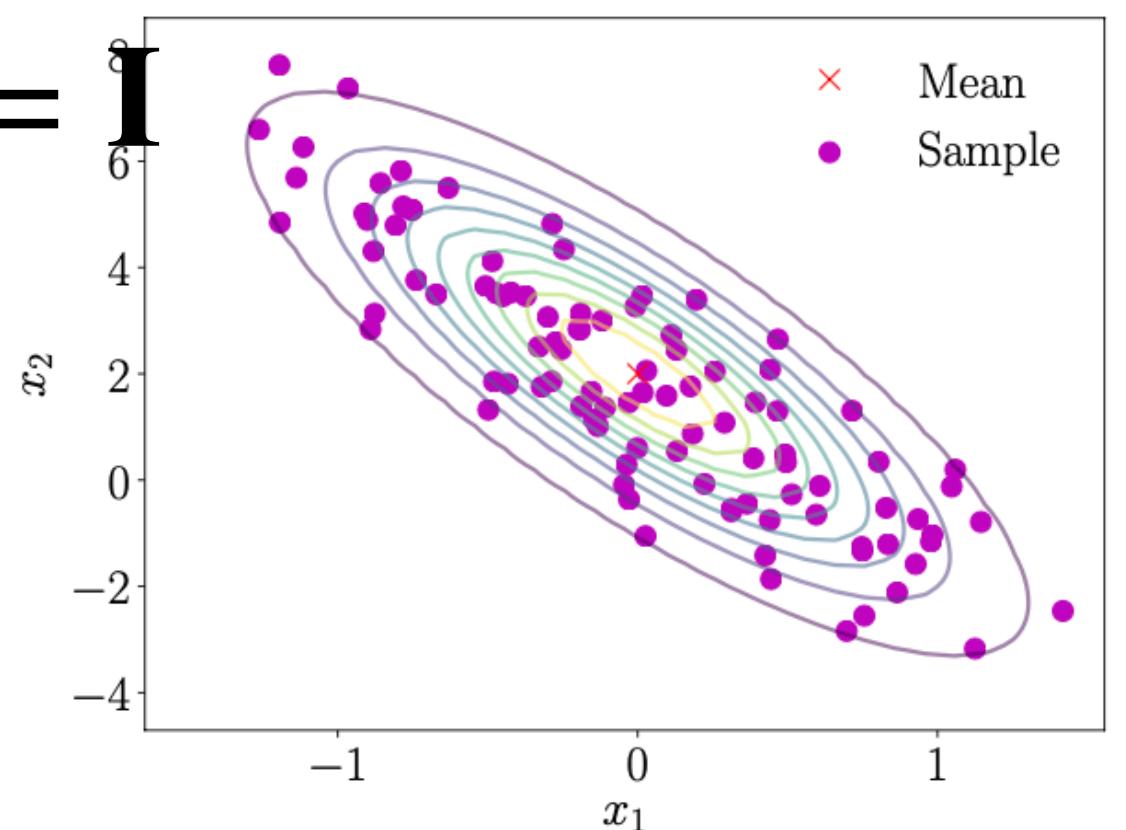
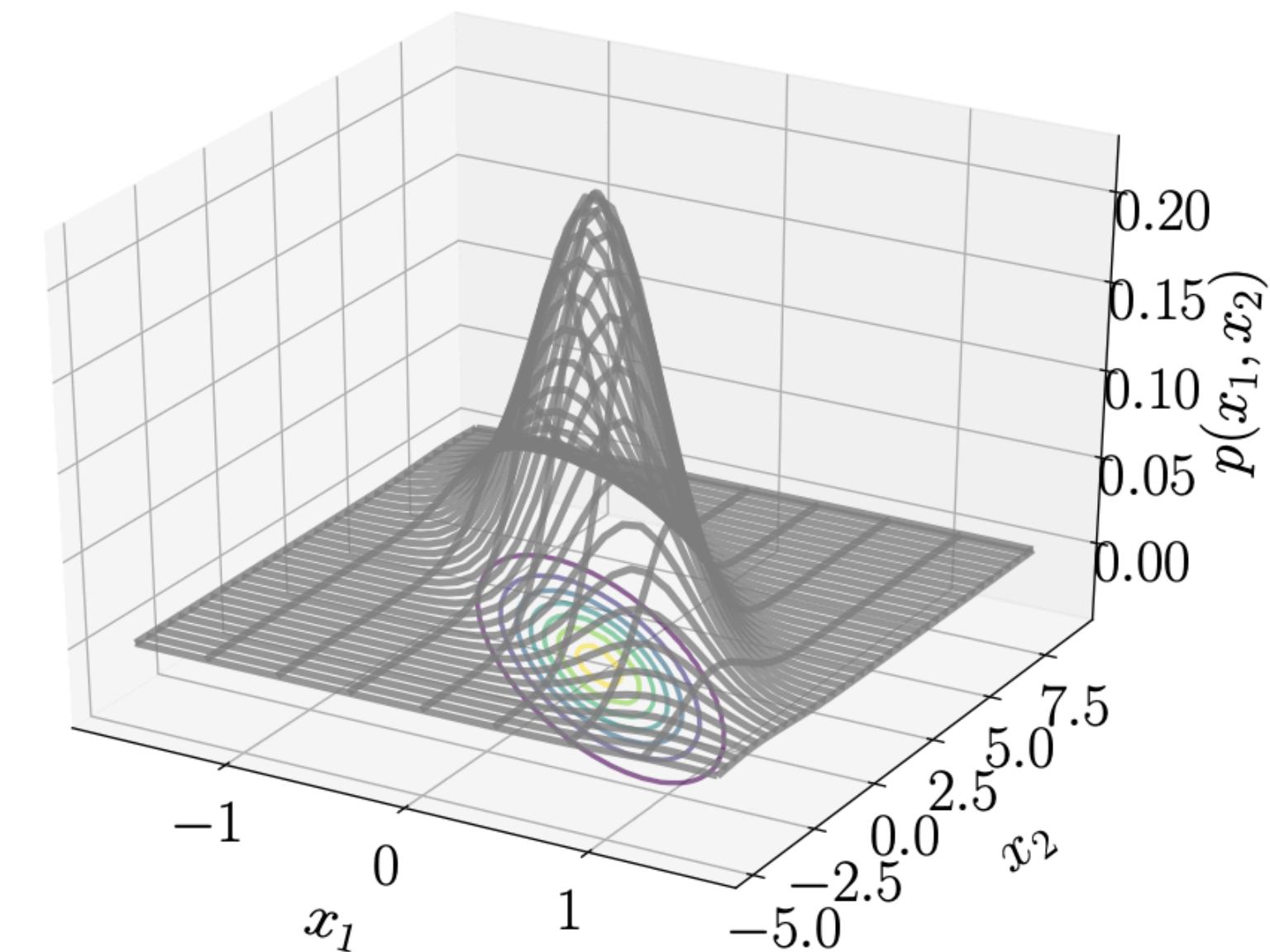
$$\text{pdf: } p(\mathbf{x} | \mu, \Sigma) = (2\pi)^{-\frac{D}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right)$$

We often write:  $p(\mathbf{x} | \mu, \Sigma) = \mathcal{N}(\mathbf{x}; \mu, \Sigma)$  or  $\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)$

Standard multivariate **Gaussian** or **normal** distribution:  $\mu = \mathbf{0}$  and  $\Sigma = \mathbf{I}$

$$\text{When } \Sigma \text{ is a diagonal matrix, } p(\mathbf{x} | \mu, \Sigma) = \prod_{d=1}^D p(x_d | \mu_d, \Sigma_{d,d})$$

Diagonal Gaussian, dimensions are independent



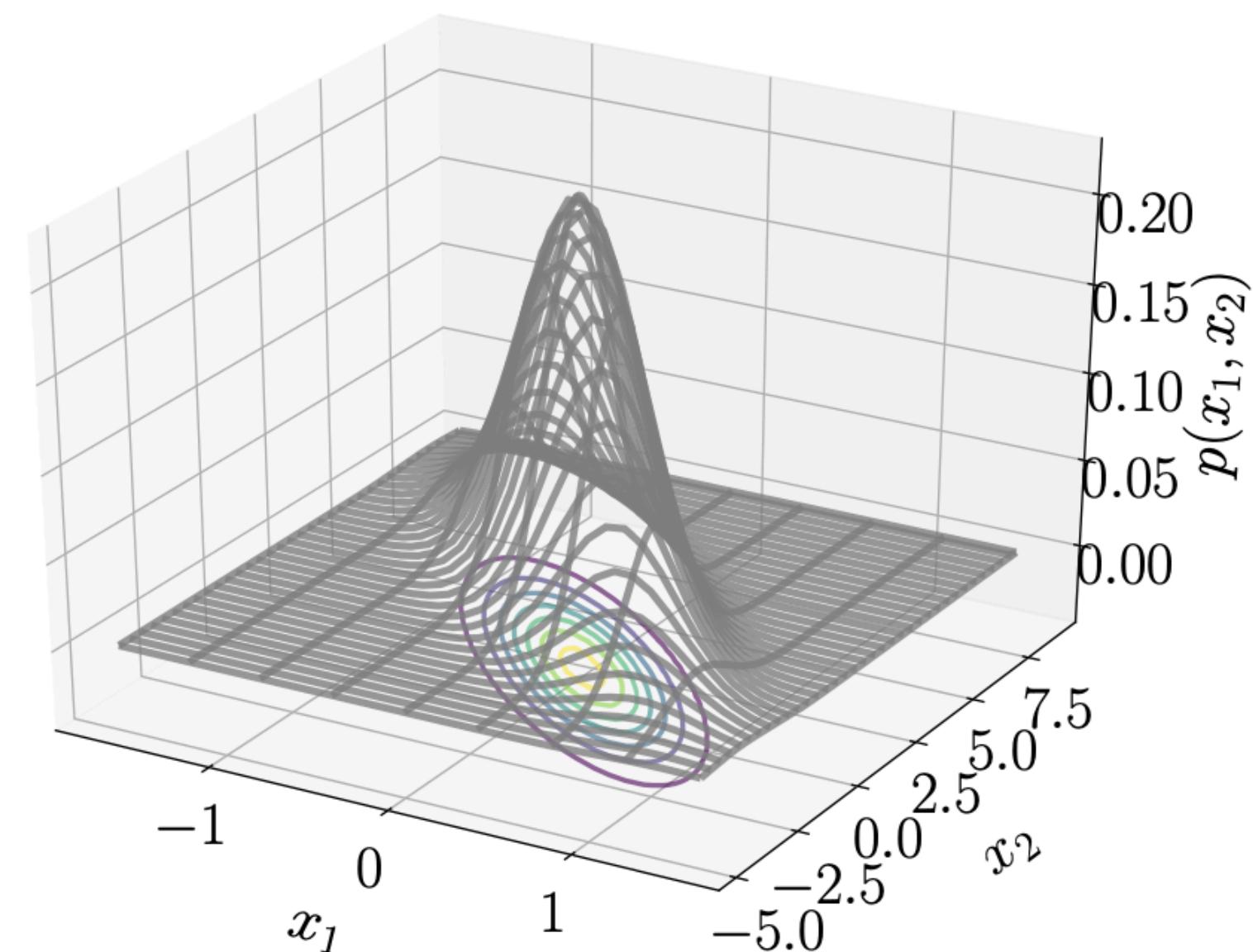
(b) Multivariate (two-dimensional) Gaussian, viewed from top. The red cross shows the mean and the colored lines show the contour lines of the density.

# Recap: multivariate Gaussian distributions (W6)

A D-dimensional multivariate **Gaussian** or **normal** distribution is characterised by a *mean* vector  $\mu$  and a *covariance matrix*  $\Sigma$ , with

$$\text{pdf: } p(\mathbf{x} | \mu, \Sigma) = (2\pi)^{-\frac{D}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right)$$

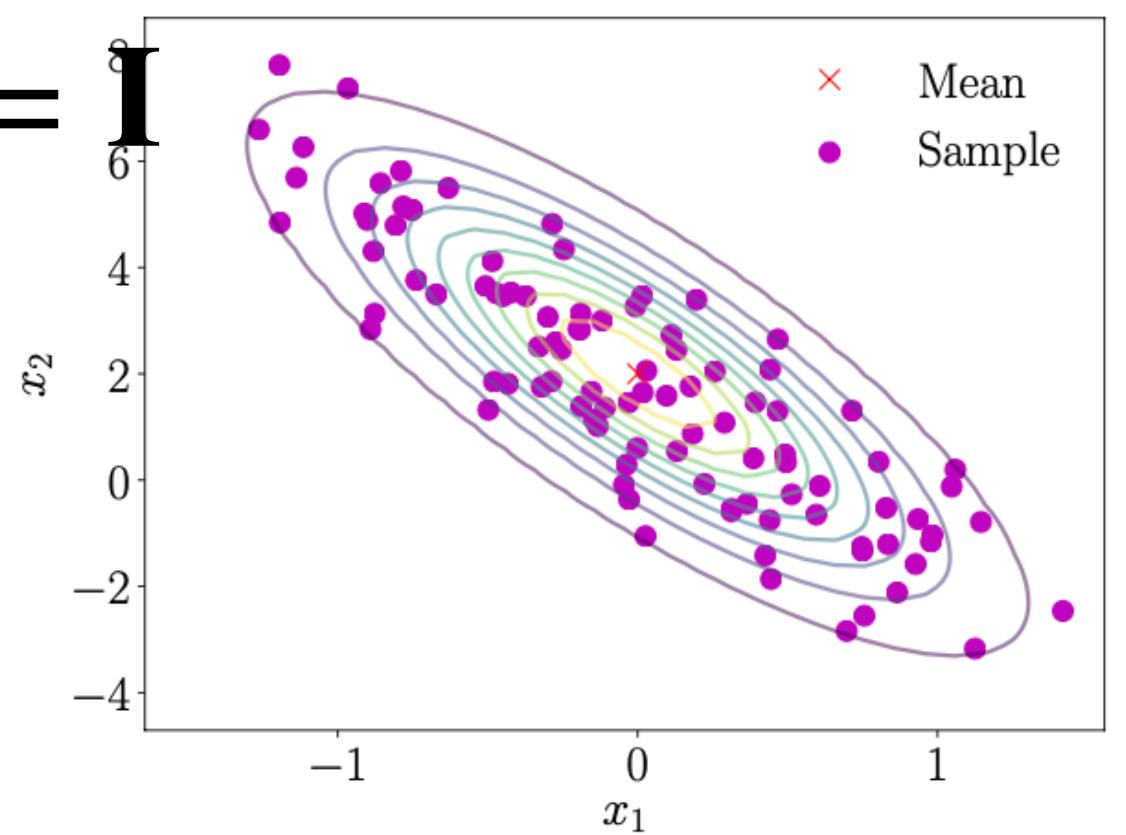
We often write:  $p(\mathbf{x} | \mu, \Sigma) = \mathcal{N}(\mathbf{x}; \mu, \Sigma)$  or  $\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)$



Standard multivariate **Gaussian** or **normal** distribution:  $\mu = \mathbf{0}$  and  $\Sigma = \mathbf{I}$

$$\text{When } \Sigma \text{ is a diagonal matrix, } p(\mathbf{x} | \mu, \Sigma) = \prod_{d=1}^D p(x_d | \mu_d, \Sigma_{d,d})$$

Diagonal Gaussian, dimensions are independent



Alternative parameterisation,  $\eta_1 = \Sigma^{-1}\mu$  and  $\eta_2 = \Sigma^{-1}$  [precision]  
Easier for multiplication/division and identifying conditional independence

(b) Multivariate (two-dimensional) Gaussian, viewed from top. The red cross shows the mean and the colored lines show the contour lines of the density.

# (Gaussian) mixture models: from one to multiple densities

# (Gaussian) mixture models: from one to multiple densities

*Mixture models* can be used to describe a distribution  $p(x)$  by a convex combination of  $K$  simple (base) distributions

$$p(x) = \sum_{k=1}^K \pi_k p_k(x) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

Why?

The components  $p_k(x)$  are, typically, members of a family of basic distributions, e.g., Gaussians, Bernoullis, or Gammas, and  $\pi_k$ 's are *mixture weights*.

# (Gaussian) mixture models: from one to multiple densities

*Mixture models* can be used to describe a distribution  $p(x)$  by a convex combination of  $K$  simple (base) distributions

$$p(x) = \sum_{k=1}^K \pi_k p_k(x) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

Why?

The components  $p_k(x)$  are, typically, members of a family of basic distributions, e.g., Gaussians, Bernoullis, or Gammas, and  $\pi_k$ 's are *mixture weights*.

When the base distributions are Gaussian distributions  $p_k(x) = \mathcal{N}(x; \mu_k, \Sigma_k)$ , we have a **Gaussian mixture model (GMM)**:

$$p_\theta(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

$\theta = \{\pi_k, \mu_k, \Sigma_k \text{ for } k = 1, 2, \dots, K\}$  are the parameters of the GMM.

# Gaussian mixture models: check your understanding

**Model:** 
$$p_{\theta}(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

where we defined  $\theta = \{\pi_k, \mu_k, \Sigma_k \text{ for } k = 1, 2, \dots, K\}$  as the parameters of the GMM.

- Given a dataset generated by a mixture of 3 Gaussians, when we randomly sample a data point, it has the probability of 1/3 belonging to each Gaussian.
- A GMM is a linear combination of several Gaussian distributions.
- In GMMs,  $K$  (number of Gaussians) is a hyper-parameter.
- If a dataset is not generated by Gaussian distributions, it cannot be modelled by a GMM.

# Gaussian mixture models: check your understanding

**Model:** 
$$p_{\theta}(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

where we defined  $\theta = \{\pi_k, \mu_k, \Sigma_k \text{ for } k = 1, 2, \dots, K\}$  as the parameters of the GMM.

- Given a dataset generated by a mixture of 3 Gaussians, when we randomly sample a data point, it has the probability of 1/3 belonging to each Gaussian. No
- A GMM is a linear combination of several Gaussian distributions.
- In GMMs,  $K$  (number of Gaussians) is a hyper-parameter.
- If a dataset is not generated by Gaussian distributions, it cannot be modelled by a GMM.

# Gaussian mixture models: check your understanding

**Model:** 
$$p_{\theta}(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

where we defined  $\theta = \{\pi_k, \mu_k, \Sigma_k \text{ for } k = 1, 2, \dots, K\}$  as the parameters of the GMM.

- Given a dataset generated by a mixture of 3 Gaussians, when we randomly sample a data point, it has the probability of 1/3 belonging to each Gaussian. No
- A GMM is a linear combination of several Gaussian distributions. Yes
- In GMMs, K (number of Gaussians) is a hyper-parameter.
- If a dataset is not generated by Gaussian distributions, it cannot be modelled by a GMM.

# Gaussian mixture models: check your understanding

**Model:** 
$$p_{\theta}(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

where we defined  $\theta = \{\pi_k, \mu_k, \Sigma_k \text{ for } k = 1, 2, \dots, K\}$  as the parameters of the GMM.

- Given a dataset generated by a mixture of 3 Gaussians, when we randomly sample a data point, it has the probability of 1/3 belonging to each Gaussian. No
- A GMM is a linear combination of several Gaussian distributions. Yes
- In GMMs, K (number of Gaussians) is a hyper-parameter. Yes
- If a dataset is not generated by Gaussian distributions, it cannot be modelled by a GMM.

# Gaussian mixture models: check your understanding

**Model:** 
$$p_{\theta}(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

where we defined  $\theta = \{\pi_k, \mu_k, \Sigma_k \text{ for } k = 1, 2, \dots, K\}$  as the parameters of the GMM.

- Given a dataset generated by a mixture of 3 Gaussians, when we randomly sample a data point, it has the probability of 1/3 belonging to each Gaussian. No
- A GMM is a linear combination of several Gaussian distributions. Yes
- In GMMs, K (number of Gaussians) is a hyper-parameter. Yes
- If a dataset is not generated by Gaussian distributions, it cannot be modelled by a GMM. No

# Gaussian mixture models: how to find $\theta$

# Gaussian mixture models: how to find $\theta$

**Model:** 
$$p_{\theta}(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

where we defined  $\theta = \{\pi_k, \mu_k, \Sigma_k \text{ for } k = 1, 2, \dots, K\}$  as the parameters of the GMM.

# Gaussian mixture models: how to find $\theta$

**Model:** 
$$p_{\theta}(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

where we defined  $\theta = \{\pi_k, \mu_k, \Sigma_k \text{ for } k = 1, 2, \dots, K\}$  as the parameters of the GMM.

**Data:**  $\mathbf{X} = \{x_1, x_2, \dots, x_N\}$ , where  $x_n$ 's are drawn i.i.d. from an unknown distribution  $p(x)$ .

# Gaussian mixture models: how to find $\theta$

**Model:** 
$$p_{\theta}(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

where we defined  $\theta = \{\pi_k, \mu_k, \Sigma_k \text{ for } k = 1, 2, \dots, K\}$  as the parameters of the GMM.

**Data:**  $\mathbf{X} = \{x_1, x_2, \dots, x_N\}$ , where  $x_n$ 's are drawn i.i.d. from an unknown distribution  $p(x)$ .

**Objective:** find a *good* approximation/representation of this unknown density  $p(x)$  by means of a GMM with  $K$  components, that is finding  $\theta$ .

# Gaussian mixture models: how to find $\theta$

**Model:** 
$$p_{\theta}(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

where we defined  $\theta = \{\pi_k, \mu_k, \Sigma_k \text{ for } k = 1, 2, \dots, K\}$  as the parameters of the GMM.

**Data:**  $\mathbf{X} = \{x_1, x_2, \dots, x_N\}$ , where  $x_n$ 's are drawn i.i.d. from an unknown distribution  $p(x)$ .

**Objective:** find a *good* approximation/representation of this unknown density  $p(x)$  by means of a GMM with  $K$  components, that is finding  $\theta$ .

$$P(\theta | \mathcal{D}) = \frac{P(\theta, \mathcal{D})}{P(\mathcal{D})} = \frac{P(\theta) P(\mathcal{D} | \theta)}{P(\mathcal{D})}$$

- Maximum likelihood (ML/MLE),  $\operatorname{argmax}_{\theta} p(\mathcal{D} | \theta)$
- Maximum a-posteriori (MAP),  $\operatorname{argmax}_{\theta} p(\mathcal{D} | \theta)p(\theta)$
- Exact/approximate Bayesian inference

Weeks 6 and 7

# Gaussian mixture models: how to find $\theta$

**Model:** 
$$p_{\theta}(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

where we defined  $\theta = \{\pi_k, \mu_k, \Sigma_k \text{ for } k = 1, 2, \dots, K\}$  as the parameters of the GMM.

**Data:**  $\mathbf{X} = \{x_1, x_2, \dots, x_N\}$ , where  $x_n$ 's are drawn i.i.d. from an unknown distribution  $p(x)$ .

**Objective:** find a *good* approximation/representation of this unknown density  $p(x)$  by means of a GMM with  $K$  components, that is finding  $\theta$ .

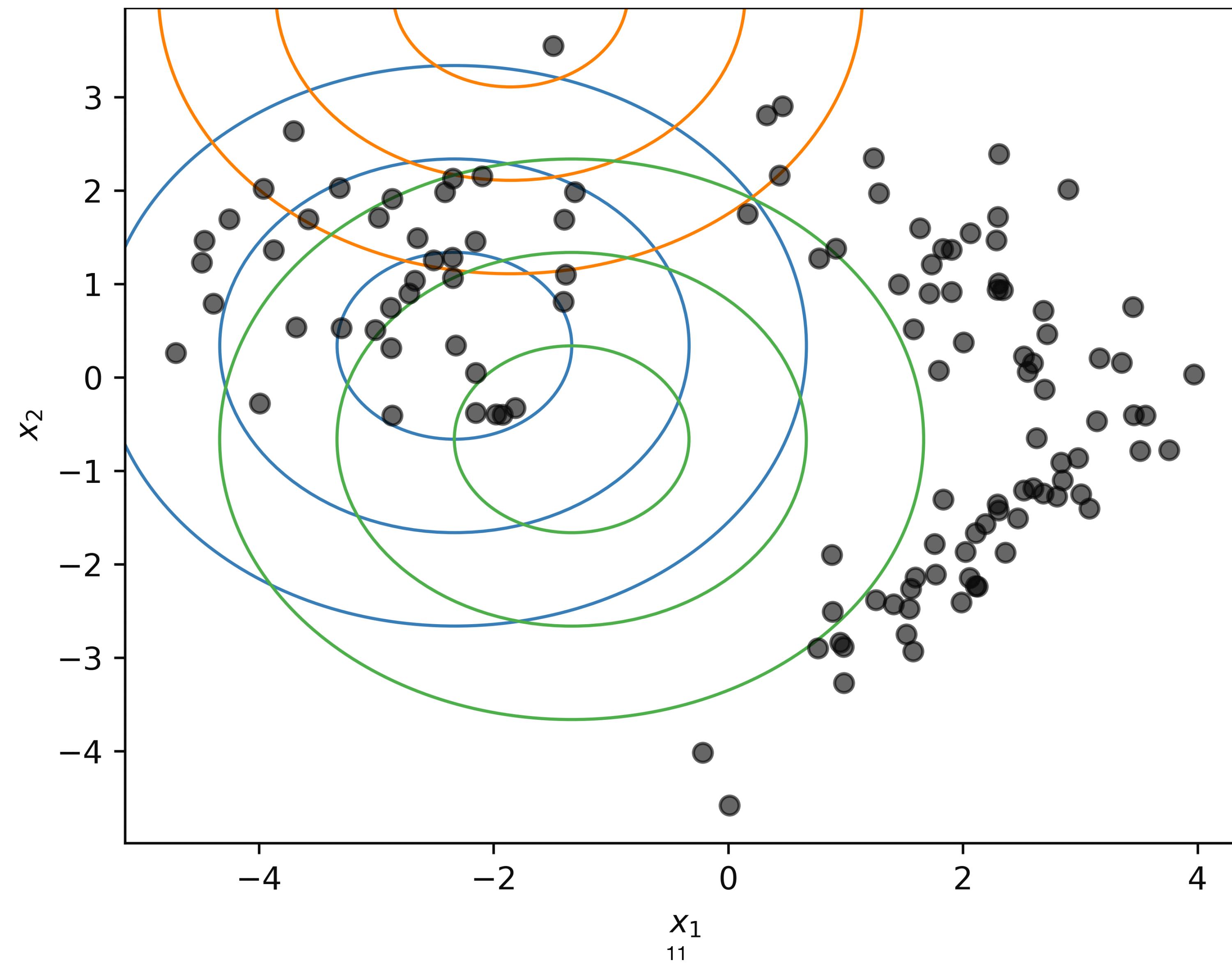
This week

$$P(\theta | \mathcal{D}) = \frac{P(\theta, \mathcal{D})}{P(\mathcal{D})} = \frac{P(\theta) P(\mathcal{D} | \theta)}{P(\mathcal{D})}$$

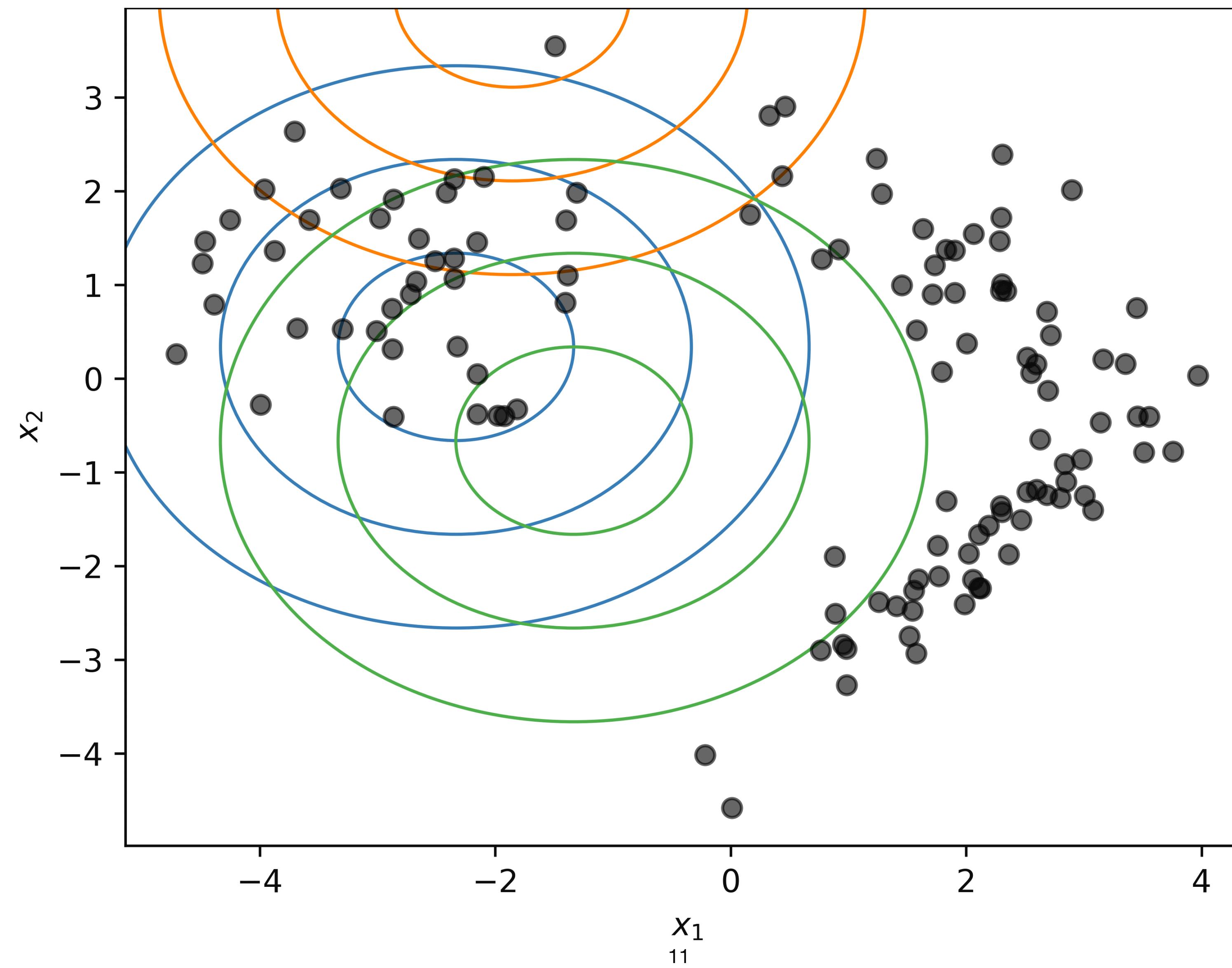
- Maximum likelihood (ML/MLE),  $\operatorname{argmax}_{\theta} p(\mathcal{D} | \theta)$
- Maximum a-posteriori (MAP),  $\operatorname{argmax}_{\theta} p(\mathcal{D} | \theta)p(\theta)$
- Exact/approximate Bayesian inference

Weeks 6 and 7

# Maximum likelihood - Sneak peek



# Maximum likelihood - Sneak peek



# GMMs - likelihood

# GMMs - likelihood

The likelihood, i.e., the predictive distribution of the training data given the parameters:

$$p(X|\theta) = \prod_{n=1}^N p(x_n|\theta) = \prod_{n=1}^N \underbrace{\sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}_{\text{Each individual likelihood term is a Gaussian mixture density}}$$

Factorise due to iid assumption

Each individual likelihood term is a Gaussian mixture density

# GMMs - likelihood

The likelihood, i.e., the predictive distribution of the training data given the parameters:

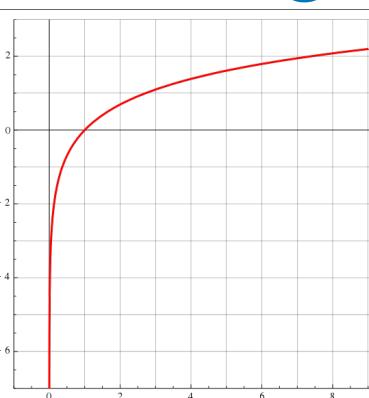
$$p(X|\theta) = \prod_{n=1}^N p(x_n|\theta) = \prod_{n=1}^N \underbrace{\sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}$$

Factorise due to iid assumption

Each individual likelihood term is a Gaussian mixture density

**Goal:** find maximise likelihood,  $\operatorname{argmax}_{\theta} p(X|\theta)$ , which is equivalent to  $\operatorname{argmax}_{\theta} \log p(X|\theta)$

as  $\log$  is monotonically increasing



# GMMs - likelihood

The likelihood, i.e., the predictive distribution of the training data given the parameters:

$$p(X|\theta) = \prod_{n=1}^N p(x_n|\theta) = \prod_{n=1}^N \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

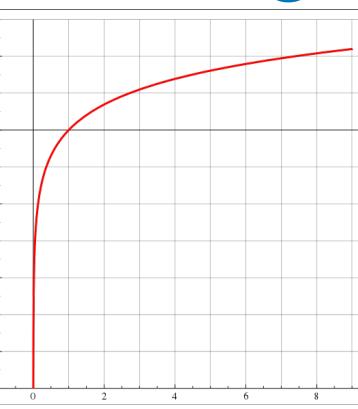
Factorise due to iid assumption

Each individual likelihood term is a Gaussian mixture density

**Goal:** find maximise likelihood,  $\operatorname{argmax}_{\theta} p(X|\theta)$ , which is equivalent to  $\operatorname{argmax}_{\theta} \log p(X|\theta)$

as log is monotonically increasing

$$\text{Log-likelihood: } L(\theta) = \log p(X|\theta) = \sum_{n=1}^N \log p(x_n|\theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$



**Strategy:** find the gradients of  $L(\theta)$  and try to set them to zero.

# GMMs - log-likelihood gradients

$$\text{Log-likelihood: } L(\theta) = \log p(X \mid \theta) = \sum_{n=1}^N \log p(x_n \mid \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

# GMMs - log-likelihood gradients

$$\text{Log-likelihood: } L(\theta) = \log p(X | \theta) = \sum_{n=1}^N \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

**Strategy:** find the gradients of  $L(\theta)$  and try to set them to zero.

$$\frac{\partial}{\partial \mu_k} L = \sum_{n=1}^N \frac{\partial \log p(x_n | \theta)}{\partial \mu_k} = \mathbf{0}^\top \quad \frac{\partial}{\partial \Sigma_k} L = \sum_{n=1}^N \frac{\partial \log p(x_n | \theta)}{\partial \Sigma_k} = \mathbf{0} \quad \frac{\partial}{\partial \pi_k} L = \sum_{n=1}^N \frac{\partial \log p(x_n | \theta)}{\partial \pi_k} = 0$$

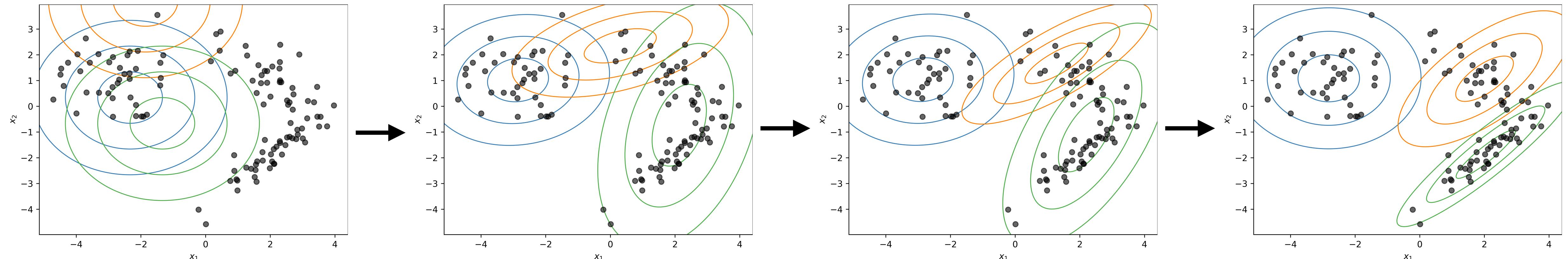
# GMMs - log-likelihood gradients

$$\text{Log-likelihood: } L(\theta) = \log p(X | \theta) = \sum_{n=1}^N \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

**Strategy:** find the gradients of  $L(\theta)$  and try to set them to zero.

$$\frac{\partial}{\partial \mu_k} L = \sum_{n=1}^N \frac{\partial \log p(x_n | \theta)}{\partial \mu_k} = \mathbf{0}^\top \quad \frac{\partial}{\partial \Sigma_k} L = \sum_{n=1}^N \frac{\partial \log p(x_n | \theta)}{\partial \Sigma_k} = \mathbf{0} \quad \frac{\partial}{\partial \pi_k} L = \sum_{n=1}^N \frac{\partial \log p(x_n | \theta)}{\partial \pi_k} = 0$$

Unlike linear regression, the optimal parameters are not analytically tractable. We resort to an iterative optimisation procedure.



Initialisation

After 50 iterations

# GMMs - iterative procedure

Initialise  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

Repeat until convergence:

1. Evaluate responsibilities  $r_{nk}$  for every data point  $x_n$  using current parameters:

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}$$

2. Re-estimate parameters  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$  using the current responsibilities  $r_{nk}$ :

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n, \quad N_k = \sum_{n=1}^N r_{nk}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^\top$$

$$\pi_k = \frac{N_k}{N}$$

# GMMs - responsibility definition

# GMMs - responsibility definition

We define **responsibility**  $r_{nk}$  for every data point  $x_n$ :

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}{\sum_j \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}$$

# GMMs - responsibility definition

We define **responsibility**  $r_{nk}$  for every data point  $x_n$ :

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}{\sum_j \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}$$

proportional to the likelihood of the  $k$ -th mixture component given the data point  $x_n$

# GMMs - responsibility definition

We define **responsibility**  $r_{nk}$  for every data point  $x_n$ :

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}{\sum_j \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}$$

proportional to the likelihood of the  $k$ -th mixture component given the data point  $x_n$

Note that  $\mathbf{r}_n = [r_{n1}, r_{n2}, \dots, r_{nK}]^\top \in \mathbb{R}^D$  is a probability vector, i.e.,  $\sum_k r_{nk} = 1$  with  $r_{nk} \geq 0$

This probability vector distributes probability mass among the  $K$  mixture components, and we can think of  $\mathbf{r}_n$  as a “soft assignment” of  $x_n$  to the  $K$  mixture components.

# GMMs - responsibility definition

We define **responsibility**  $r_{nk}$  for every data point  $x_n$ :

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}{\sum_j \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}$$

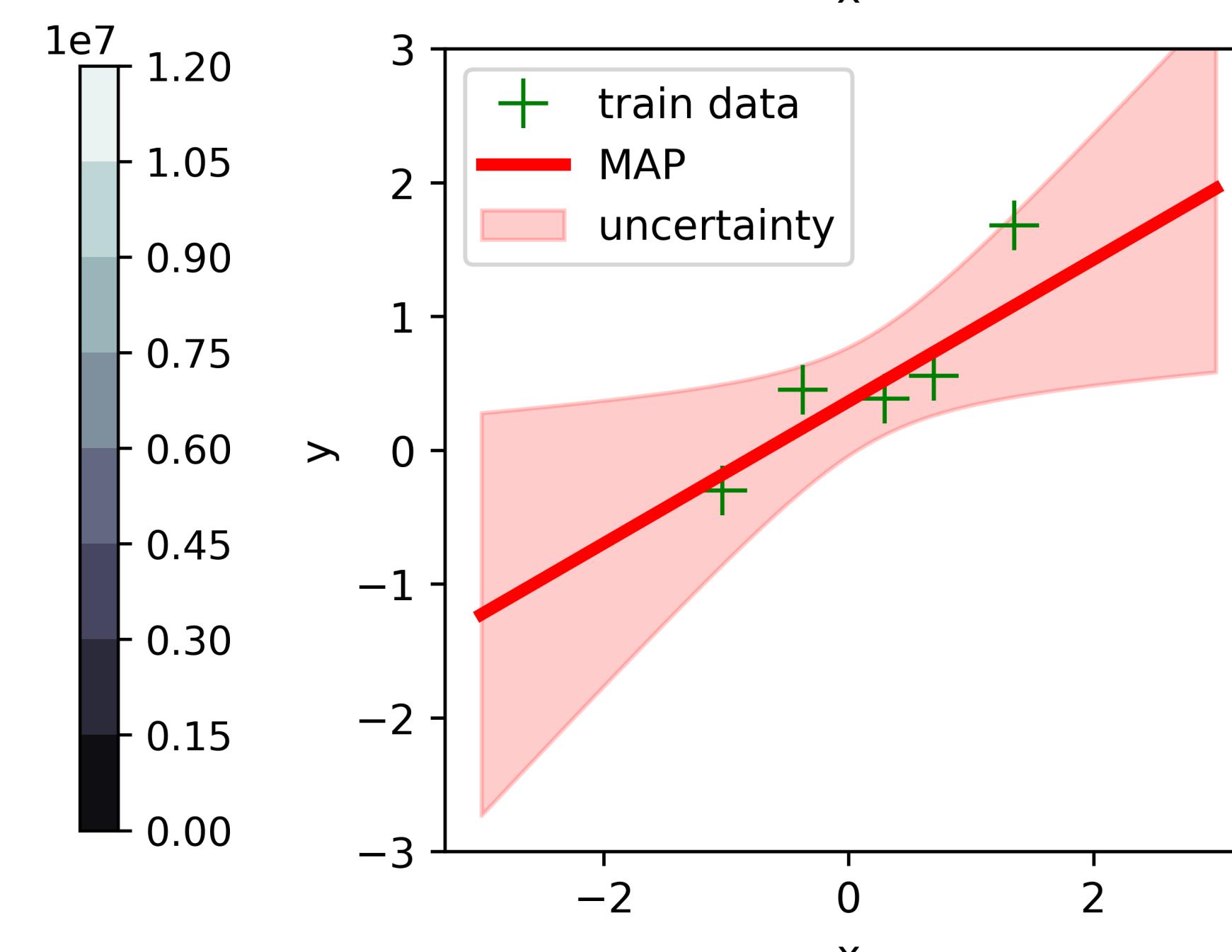
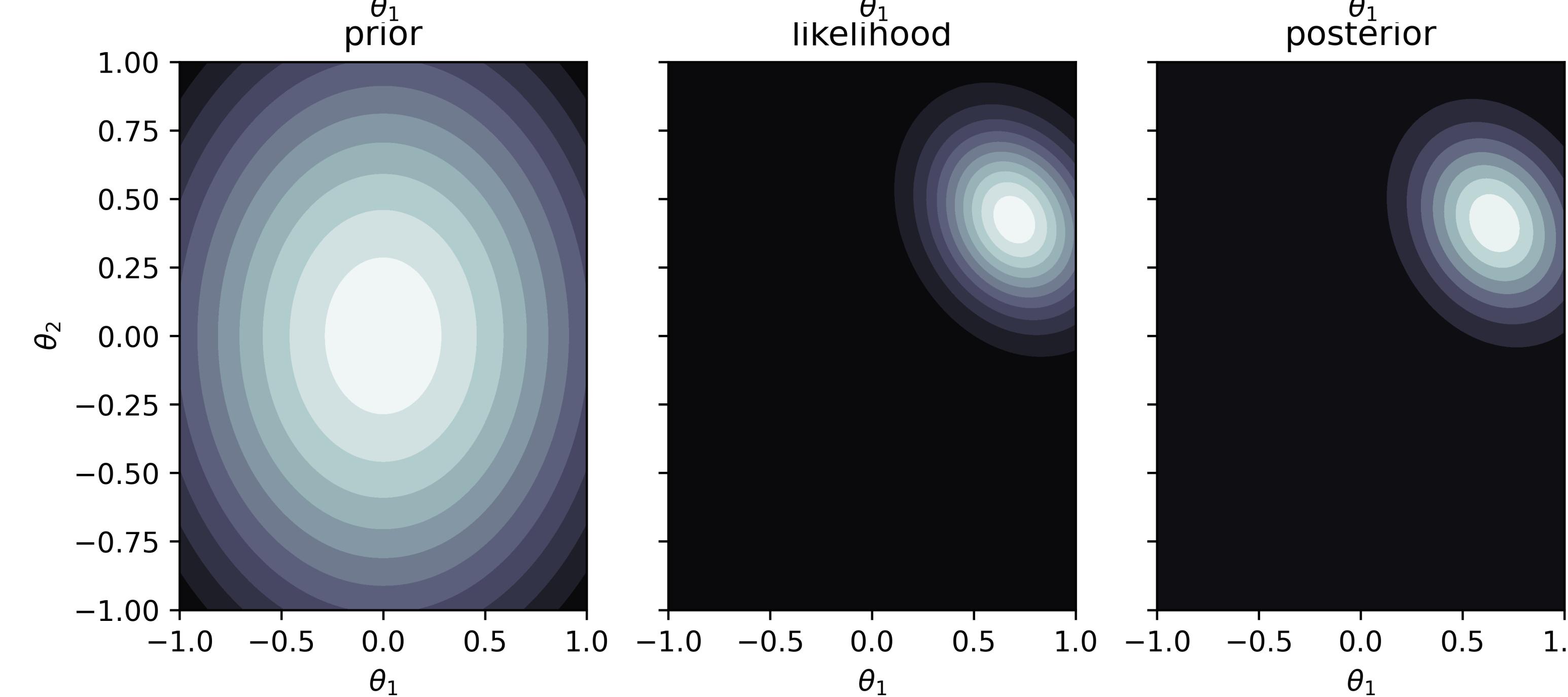
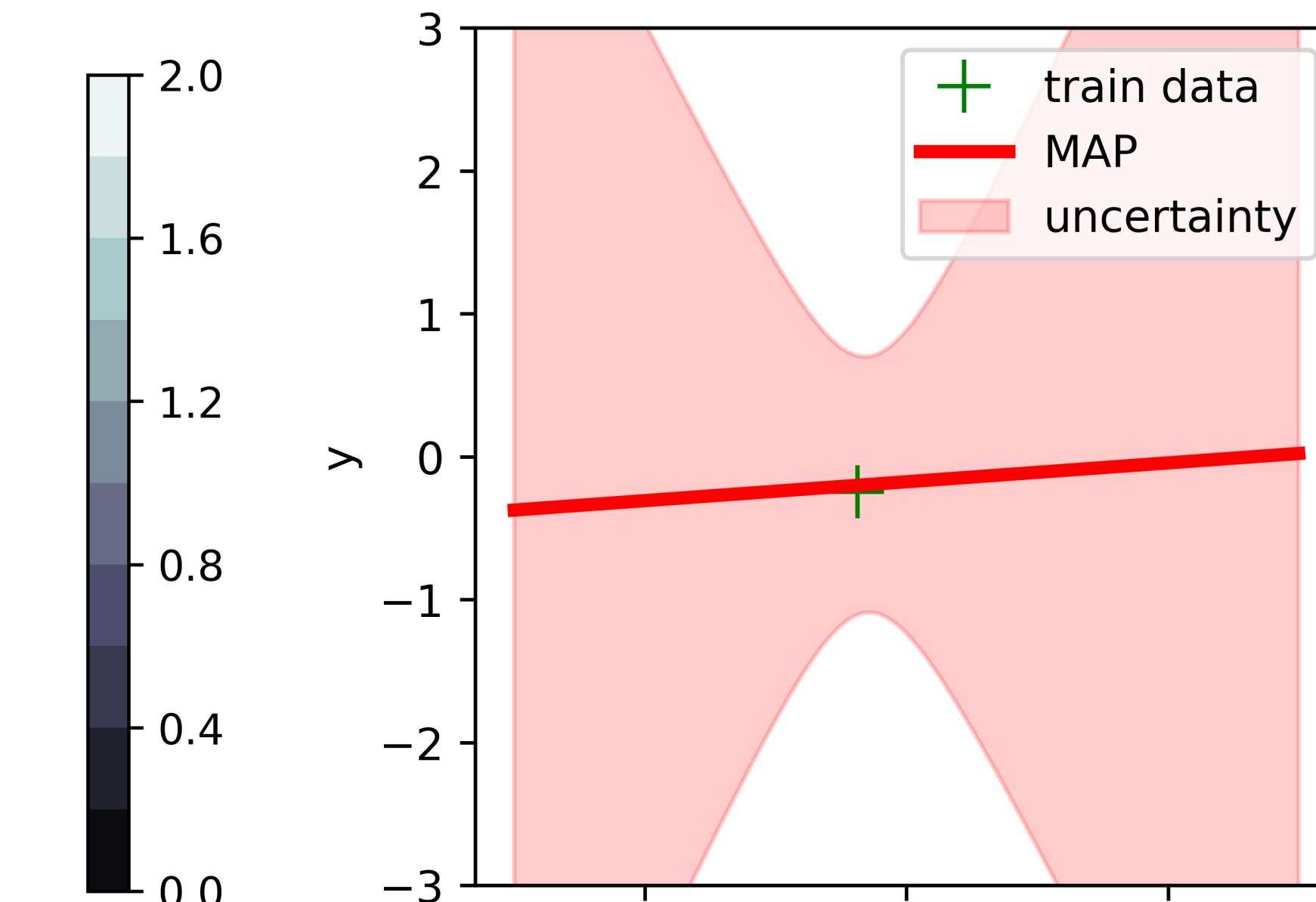
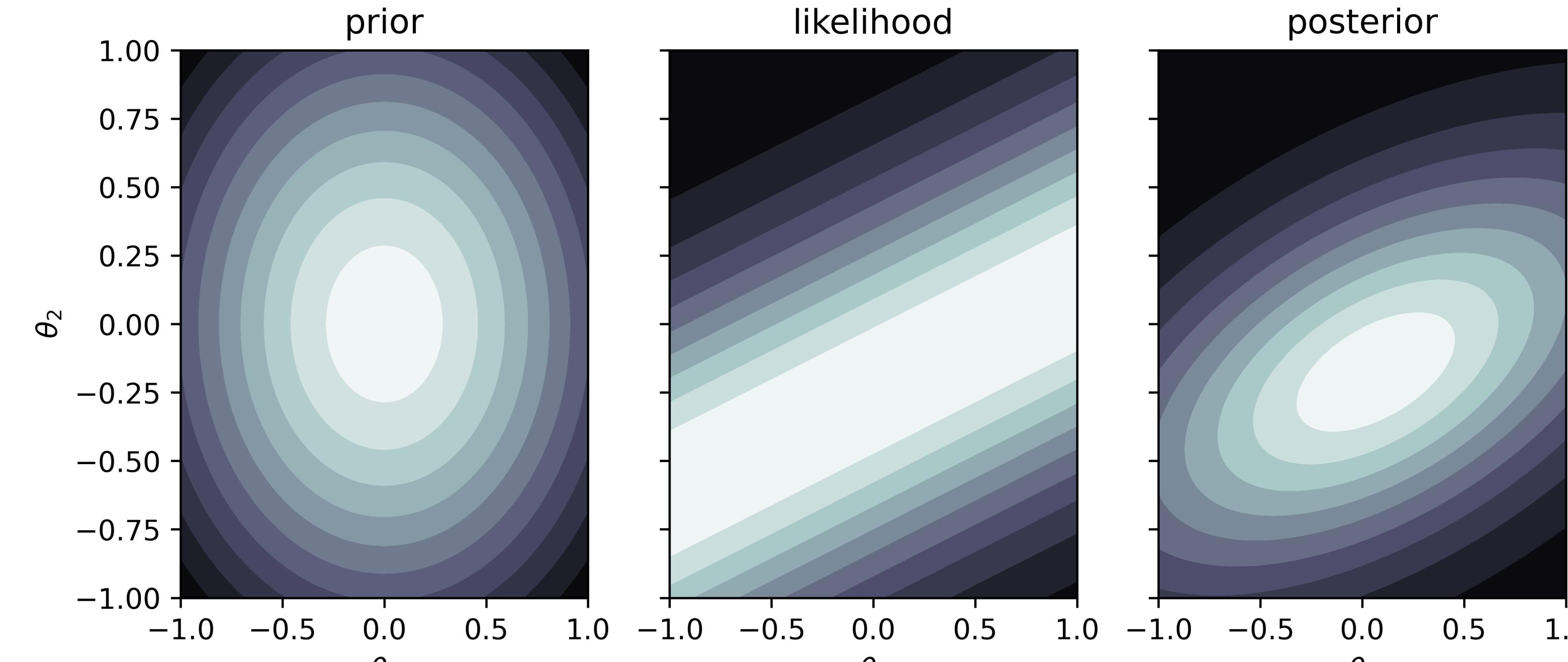
proportional to the likelihood of the  $k$ -th mixture component given the data point  $x_n$

Note that  $\mathbf{r}_n = [r_{n1}, r_{n2}, \dots, r_{nK}]^\top \in \mathbb{R}^D$  is a probability vector, i.e.,  $\sum_k r_{nk} = 1$  with  $r_{nk} \geq 0$

This probability vector distributes probability mass among the  $K$  mixture components, and we can think of  $\mathbf{r}_n$  as a “soft assignment” of  $x_n$  to the  $K$  mixture components.

Let's now derive the parameter updates in **step 2!**

# More on Bayesian linear regression



# Overview

Last lecture:

1. **Model:** Gaussian mixture models
2. **Learning:** parameter learning using *maximum likelihood*
3. **Iterative** procedure, derive mean update

**Initialise**  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

**Repeat** until convergence:

1. **Evaluate responsibilities**  $r_{nk}$  for every data point  $x_n$  using current parameters:

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}$$

2. **Re-estimate parameters**  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$  using the current responsibilities  $r_{nk}$ :

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n, \quad N_k = \sum_{n=1}^N r_{nk}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^\top$$

$$\pi_k = \frac{N_k}{N}$$

This lecture:

1. Derive covariance and mixture weight updates
2. Issues with maximum likelihood and workarounds
3. K-means as a special case
4. Theoretical foundation: EM algorithm (1 is E step, 2 is M step)

# Recap: Gaussian mixture models and how to find $\theta$

**Model:** 
$$p_{\theta}(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

where we defined  $\theta = \{\pi_k, \mu_k, \Sigma_k \text{ for } k = 1, 2, \dots, K\}$  as the parameters of the GMM.

**Data:**  $\mathbf{X} = \{x_1, x_2, \dots, x_N\}$ , where  $x_n$ 's are drawn i.i.d. from an unknown distribution  $p(x)$ .

**Objective:** find a *good* approximation/representation of this unknown density  $p(x)$  by means of a GMM with  $K$  components, that is finding  $\theta$ .

This week

$$P(\theta | \mathcal{D}) = \frac{P(\theta, \mathcal{D})}{P(\mathcal{D})} = \frac{P(\theta) P(\mathcal{D} | \theta)}{P(\mathcal{D})}$$

- Maximum likelihood (ML/MLE),  $\operatorname{argmax}_{\theta} p(\mathcal{D} | \theta)$
- Maximum a-posteriori (MAP),  $\operatorname{argmax}_{\theta} p(\mathcal{D} | \theta)p(\theta)$
- Exact/approximate Bayesian inference

Weeks 6 and 7

# Recap: GMMs - iterative procedure

Initialise  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

Repeat until convergence:

1. Evaluate responsibilities  $r_{nk}$  for every data point  $x_n$  using current parameters:

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}$$

2. Re-estimate parameters  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$  using the current responsibilities  $r_{nk}$ :

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n, \quad N_k = \sum_{n=1}^N r_{nk}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^\top$$

$$\pi_k = \frac{N_k}{N}$$

# Recap: Step 2, mean update

$$L(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k) \quad \text{Log likelihood}$$

# Recap: Step 2, mean update

$$L(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Log likelihood

$$\frac{\partial L(\theta)}{\partial \mu_k} = \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \frac{\partial}{\partial \mu_k} \sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)$$

gradient of log likelihood wrt mean

# Recap: Step 2, mean update

$$L(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Log likelihood

$$\frac{\partial L(\theta)}{\partial \mu_k} = \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \frac{\partial}{\partial \mu_k} \sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)$$

$$= \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \pi_k \frac{\partial}{\partial \mu_k} \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

gradient of log likelihood wrt mean

Ignore terms that do not depend on  $k$

# Recap: Step 2, mean update

$$L(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Log likelihood

$$\frac{\partial L(\theta)}{\partial \mu_k} = \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \frac{\partial}{\partial \mu_k} \sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)$$

gradient of log likelihood wrt mean

$$= \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \pi_k \frac{\partial}{\partial \mu_k} \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Ignore terms that do not depend on k

$$= \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k) (x_n - \mu_k)^\top \Sigma_k^{-1}$$

Gradient of Gaussian wrt mean

# Recap: Step 2, mean update

$$L(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Log likelihood

$$\frac{\partial L(\theta)}{\partial \mu_k} = \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \frac{\partial}{\partial \mu_k} \sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)$$

gradient of log likelihood wrt mean

$$= \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \pi_k \frac{\partial}{\partial \mu_k} \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Ignore terms that do not depend on  $k$

$$= \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k) (x_n - \mu_k)^\top \Sigma_k^{-1}$$

Gradient of Gaussian wrt mean

$$= \sum_{n=1}^N r_{nk} (x_n - \mu_k)^\top \Sigma_k^{-1}$$

Rewrite using responsibility

# Recap: Step 2, mean update

$$L(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Log likelihood

$$\frac{\partial L(\theta)}{\partial \mu_k} = \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \frac{\partial}{\partial \mu_k} \sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)$$

gradient of log likelihood wrt mean

$$= \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \pi_k \frac{\partial}{\partial \mu_k} \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

$$= \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k) (x_n - \mu_k)^\top \Sigma_k^{-1}$$

$$= \sum_{n=1}^N r_{nk} (x_n - \mu_k)^\top \Sigma_k^{-1}$$

$$\frac{\partial L(\theta)}{\partial \mu_k} = \mathbf{0}^\top \rightarrow \mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}}$$

Set gradient to zero

Ignore terms that do not depend on k

Gradient of Gaussian wrt mean

Rewrite using responsibility

# Step 2, covariance update [1]

$$L(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Log likelihood

# Step 2, covariance update [1]

$$L(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Log likelihood

$$\frac{\partial L(\theta)}{\partial \Sigma_k} = \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \frac{\partial}{\partial \Sigma_k} \sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)$$

gradient of log likelihood wrt cov

$$= \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \pi_k \frac{\partial}{\partial \Sigma_k} \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Ignore terms that do not depend on k

# Step 2, covariance update [1]

$$L(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Log likelihood

$$\begin{aligned} \frac{\partial L(\theta)}{\partial \Sigma_k} &= \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \frac{\partial}{\partial \Sigma_k} \sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j) \\ &= \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \pi_k \frac{\partial}{\partial \Sigma_k} \mathcal{N}(x_n; \mu_k, \Sigma_k) \end{aligned}$$

gradient of log likelihood wrt cov

Ignore terms that do not depend on k

$$\frac{\partial}{\partial \Sigma_k} \mathcal{N}(x_n; \mu_k, \Sigma_k) = \frac{\partial}{\partial \Sigma_k} \left[ (2\pi)^{-\frac{D}{2}} |\Sigma_k|^{-\frac{1}{2}} \exp \left( -\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) \right) \right]$$

Gradient of Gaussian  
wrt cov

# Step 2, covariance update [1]

$$L(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Log likelihood

$$\frac{\partial L(\theta)}{\partial \Sigma_k} = \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \frac{\partial}{\partial \Sigma_k} \sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)$$

gradient of log likelihood wrt cov

$$= \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \pi_k \frac{\partial}{\partial \Sigma_k} \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Ignore terms that do not depend on k

$$\frac{\partial}{\partial \Sigma_k} \mathcal{N}(x_n; \mu_k, \Sigma_k) = \frac{\partial}{\partial \Sigma_k} \left[ (2\pi)^{-\frac{D}{2}} |\Sigma_k|^{-\frac{1}{2}} \exp \left( -\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) \right) \right]$$

Gradient of Gaussian  
wrt cov

$$= (2\pi)^{-\frac{D}{2}} \frac{-1}{2} |\Sigma_k|^{-\frac{1}{2}} \Sigma_k^{-1} \exp \left( -\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) \right)$$

$$+ (2\pi)^{-\frac{D}{2}} |\Sigma_k|^{-\frac{1}{2}} \exp \left( -\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) \right) \frac{1}{2} \Sigma_k^{-1} (x_n - \mu_k) (x_n - \mu_k)^T \Sigma_k^{-1}$$

# Step 2, covariance update [1]

$$L(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Log likelihood

$$\frac{\partial L(\theta)}{\partial \Sigma_k} = \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \frac{\partial}{\partial \Sigma_k} \sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)$$

gradient of log likelihood wrt cov

$$= \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \pi_k \frac{\partial}{\partial \Sigma_k} \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Ignore terms that do not depend on k

$$\frac{\partial}{\partial \Sigma_k} \mathcal{N}(x_n; \mu_k, \Sigma_k) = \frac{\partial}{\partial \Sigma_k} \left[ (2\pi)^{-\frac{D}{2}} |\Sigma_k|^{-\frac{1}{2}} \exp \left( -\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) \right) \right]$$

Gradient of Gaussian  
wrt cov

$$= (2\pi)^{-\frac{D}{2}} \frac{-1}{2} |\Sigma_k|^{-\frac{1}{2}} \Sigma_k^{-1} \exp \left( -\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) \right)$$

$$+ (2\pi)^{-\frac{D}{2}} |\Sigma_k|^{-\frac{1}{2}} \exp \left( -\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) \right) \frac{1}{2} \Sigma_k^{-1} (x_n - \mu_k) (x_n - \mu_k)^T \Sigma_k^{-1}$$

$$= -\frac{1}{2} \mathcal{N}(x_n; \mu_k, \Sigma_k) [\Sigma_k^{-1} - \Sigma_k^{-1} (x_n - \mu_k) (x_n - \mu_k)^T \Sigma_k^{-1}]$$

# Step 2, covariance update [2]

$$\frac{\partial L(\theta)}{\partial \Sigma_k} = \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \pi_k \frac{-1}{2} \mathcal{N}(x_n; \mu_k, \Sigma_k) [\Sigma_k^{-1} - \Sigma_k^{-1} (x_n - \mu_k) (x_n - \mu_k)^\top \Sigma_k^{-1}]$$

gradient of log likelihood wrt cov

# Step 2, covariance update [2]

$$\frac{\partial L(\theta)}{\partial \Sigma_k} = \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \pi_k \frac{-1}{2} \mathcal{N}(x_n; \mu_k, \Sigma_k) [\Sigma_k^{-1} - \Sigma_k^{-1} (x_n - \mu_k) (x_n - \mu_k)^\top \Sigma_k^{-1}]$$

gradient of log likelihood wrt cov

$$= -\frac{1}{2} \sum_{n=1}^N \textcolor{red}{r}_{nk} [\Sigma_k^{-1} - \Sigma_k^{-1} (x_n - \mu_k) (x_n - \mu_k)^\top \Sigma_k^{-1}]$$

Rewrite using responsibility

# Step 2, covariance update [2]

$$\frac{\partial L(\theta)}{\partial \Sigma_k} = \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \pi_k \frac{-1}{2} \mathcal{N}(x_n; \mu_k, \Sigma_k) [\Sigma_k^{-1} - \Sigma_k^{-1} (x_n - \mu_k) (x_n - \mu_k)^\top \Sigma_k^{-1}]$$

gradient of log likelihood wrt cov

$$= -\frac{1}{2} \sum_{n=1}^N \textcolor{red}{r}_{nk} [\Sigma_k^{-1} - \Sigma_k^{-1} (x_n - \mu_k) (x_n - \mu_k)^\top \Sigma_k^{-1}]$$

Rewrite using responsibility

$$\frac{\partial L(\theta)}{\partial \Sigma_k} = 0 \rightarrow \Sigma_k = \frac{\sum_n r_{nk} (x_n - \mu_k) (x_n - \mu_k)^\top}{\sum_n r_{nk}}$$

Set gradient to zero

# Step 2, covariance update [2]

$$\frac{\partial L(\theta)}{\partial \Sigma_k} = \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \pi_k \frac{-1}{2} \mathcal{N}(x_n; \mu_k, \Sigma_k) [\Sigma_k^{-1} - \Sigma_k^{-1} (x_n - \mu_k) (x_n - \mu_k)^\top \Sigma_k^{-1}]$$

gradient of log likelihood wrt cov

$$= -\frac{1}{2} \sum_{n=1}^N \textcolor{red}{r}_{nk} [\Sigma_k^{-1} - \Sigma_k^{-1} (x_n - \mu_k) (x_n - \mu_k)^\top \Sigma_k^{-1}]$$

Rewrite using responsibility

$$\frac{\partial L(\theta)}{\partial \Sigma_k} = 0 \rightarrow \Sigma_k = \frac{\sum_n r_{nk} (x_n - \mu_k) (x_n - \mu_k)^\top}{\sum_n r_{nk}}$$

Set gradient to zero

New covariance = weighted covariance of data where weights = responsibilities

# Step 2, mixture weight update

$$L(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k) \quad \text{Log likelihood}$$

# Step 2, mixture weight update

$$L(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Log likelihood

$$\hat{L}(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k) + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

Lagrangian to deal with equality constraint  $\sum_k \pi_k = 1$

Covered in optimisation and advanced ML

# Step 2, mixture weight update

$$L(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Log likelihood

$$\hat{L}(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k) + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

Lagrangian to deal with  
equality constraint

Covered in optimisation  
and advanced ML

$$\sum_k \pi_k = 1$$

$$\frac{\partial \hat{L}(\theta)}{\partial \pi_k} = \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \frac{\partial}{\partial \pi_k} \sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j) + \frac{\partial}{\partial \pi_k} \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

gradient of Lagrangian wrt cov

# Step 2, mixture weight update

$$L(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Log likelihood

$$\hat{L}(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k) + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

Lagrangian to deal with  
equality constraint

Covered in optimisation  
and advanced ML

$$\sum_k \pi_k = 1$$

gradient of Lagrangian wrt cov

$$\frac{\partial \hat{L}(\theta)}{\partial \pi_k} = \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \frac{\partial}{\partial \pi_k} \sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j) + \frac{\partial}{\partial \pi_k} \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

$$= \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \mathcal{N}(x_n; \mu_k, \Sigma_k) + \lambda$$

Ignore terms that do not depend on k

# Step 2, mixture weight update

$$L(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Log likelihood

$$\hat{L}(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k) + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

Lagrangian to deal with equality constraint

Covered in optimisation and advanced ML

$$\sum_k \pi_k = 1$$

gradient of Lagrangian wrt cov

$$\frac{\partial \hat{L}(\theta)}{\partial \pi_k} = \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \frac{\partial}{\partial \pi_k} \sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j) + \frac{\partial}{\partial \pi_k} \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

$$= \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \mathcal{N}(x_n; \mu_k, \Sigma_k) + \lambda$$

Ignore terms that do not depend on k

$$= \sum_{n=1}^N \frac{r_{nk}}{\pi_k} + \lambda$$

Rewrite using responsibility

# Step 2, mixture weight update

$$L(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Log likelihood

$$\hat{L}(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k) + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

Lagrangian to deal with equality constraint

Covered in optimisation and advanced ML

$$\sum_k \pi_k = 1$$

gradient of Lagrangian wrt cov

$$\frac{\partial \hat{L}(\theta)}{\partial \pi_k} = \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \frac{\partial}{\partial \pi_k} \sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j) + \frac{\partial}{\partial \pi_k} \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

$$= \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \mathcal{N}(x_n; \mu_k, \Sigma_k) + \lambda$$

Ignore terms that do not depend on k

$$= \sum_{n=1}^N \frac{r_{nk}}{\pi_k} + \lambda$$

Rewrite using responsibility

$$\frac{\partial \hat{L}(\theta)}{\partial \lambda} = \sum_{k=1}^K \pi_k - 1$$

gradient wrt Lagrange multiplier

# Step 2, mixture weight update

$$L(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Log likelihood

$$\hat{L}(\theta) = \sum_n \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k) + \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

Lagrangian to deal with equality constraint  $\sum_k \pi_k = 1$

$$\frac{\partial \hat{L}(\theta)}{\partial \pi_k} = \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \frac{\partial}{\partial \pi_k} \sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j) + \frac{\partial}{\partial \pi_k} \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

gradient of Lagrangian wrt cov

$$= \sum_{n=1}^N \frac{1}{p(x_n | \theta)} \mathcal{N}(x_n; \mu_k, \Sigma_k) + \lambda$$

Ignore terms that do not depend on  $k$

$$= \sum_{n=1}^N \frac{r_{nk}}{\pi_k} + \lambda$$

Rewrite using responsibility

$$\frac{\partial \hat{L}(\theta)}{\partial \lambda} = \sum_{k=1}^K \pi_k - 1$$

gradient wrt Lagrange multiplier

$$\frac{\partial \hat{L}(\theta)}{\partial \pi_k} = 0 \quad \text{and} \quad \frac{\partial \hat{L}(\theta)}{\partial \lambda} = 0 \rightarrow \pi_k = \frac{\sum_n r_{nk}}{N}$$

Set gradients to zero

Covered in optimisation and advanced ML

# Overview

Last lecture:

1. **Model:** Gaussian mixture models

2. **Learning:** parameter learning using *maximum likelihood*

3. **Iterative** procedure, derive mean update

**Initialise**  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

**Repeat** until convergence:

1. **Evaluate responsibilities**  $r_{nk}$  for every data point  $x_n$  using current parameters:

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}$$

2. **Re-estimate parameters**  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$  using the current responsibilities  $r_{nk}$ :

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n, \quad N_k = \sum_{n=1}^N r_{nk}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^\top$$

$$\pi_k = \frac{N_k}{N}$$

This lecture:

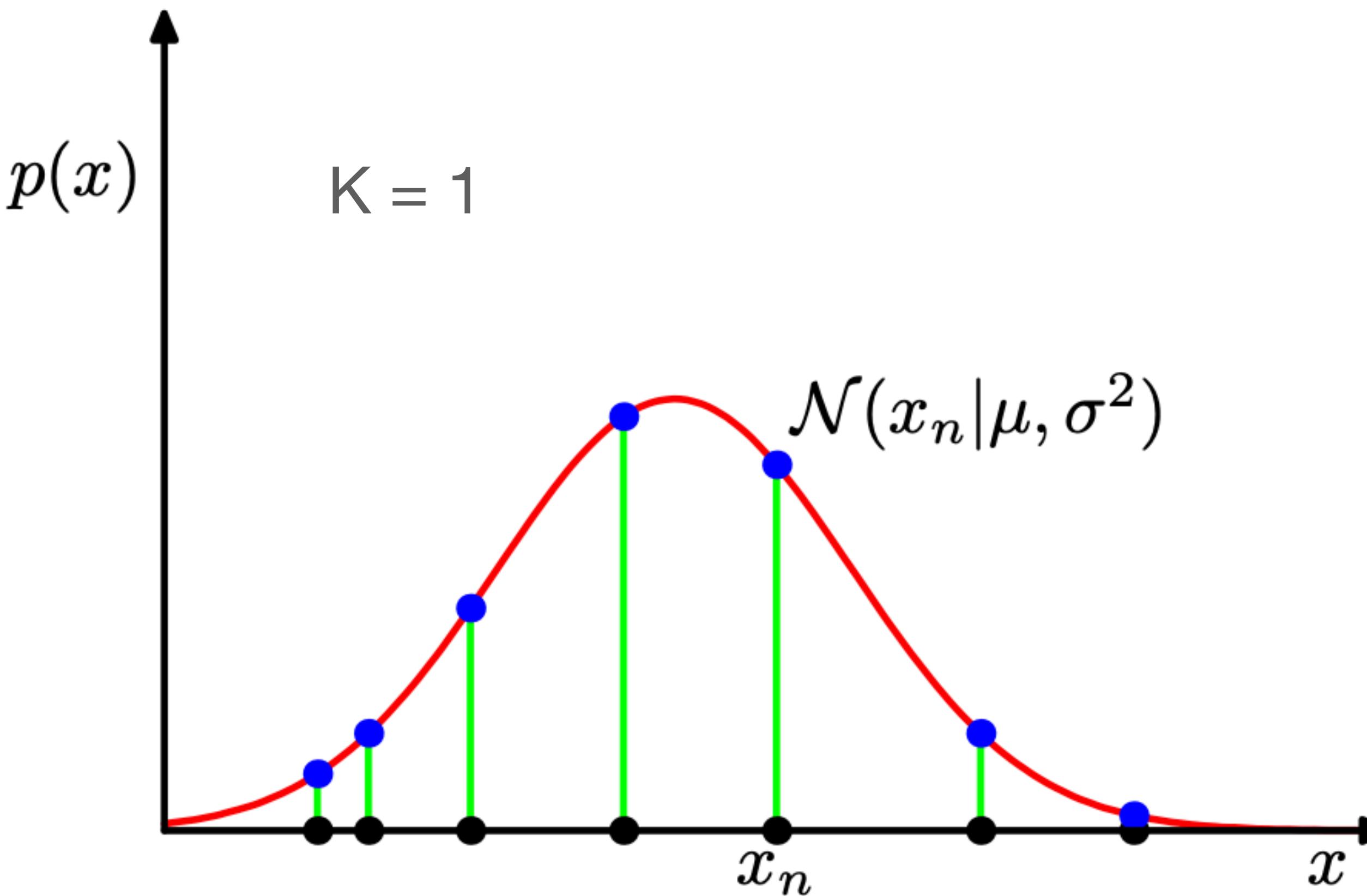
1. Derive covariance and mixture weight updates

2. Issues with maximum likelihood and workarounds

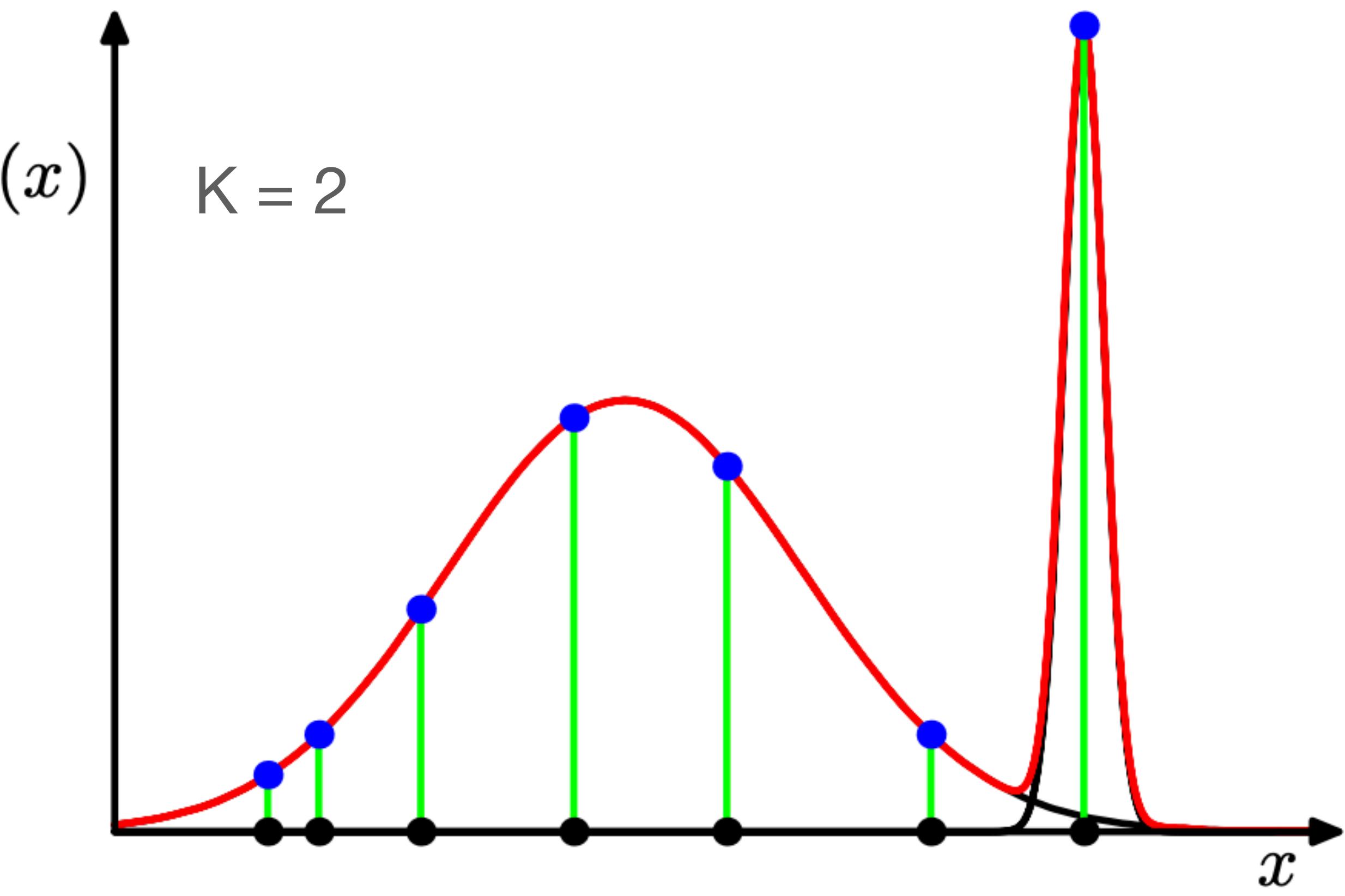
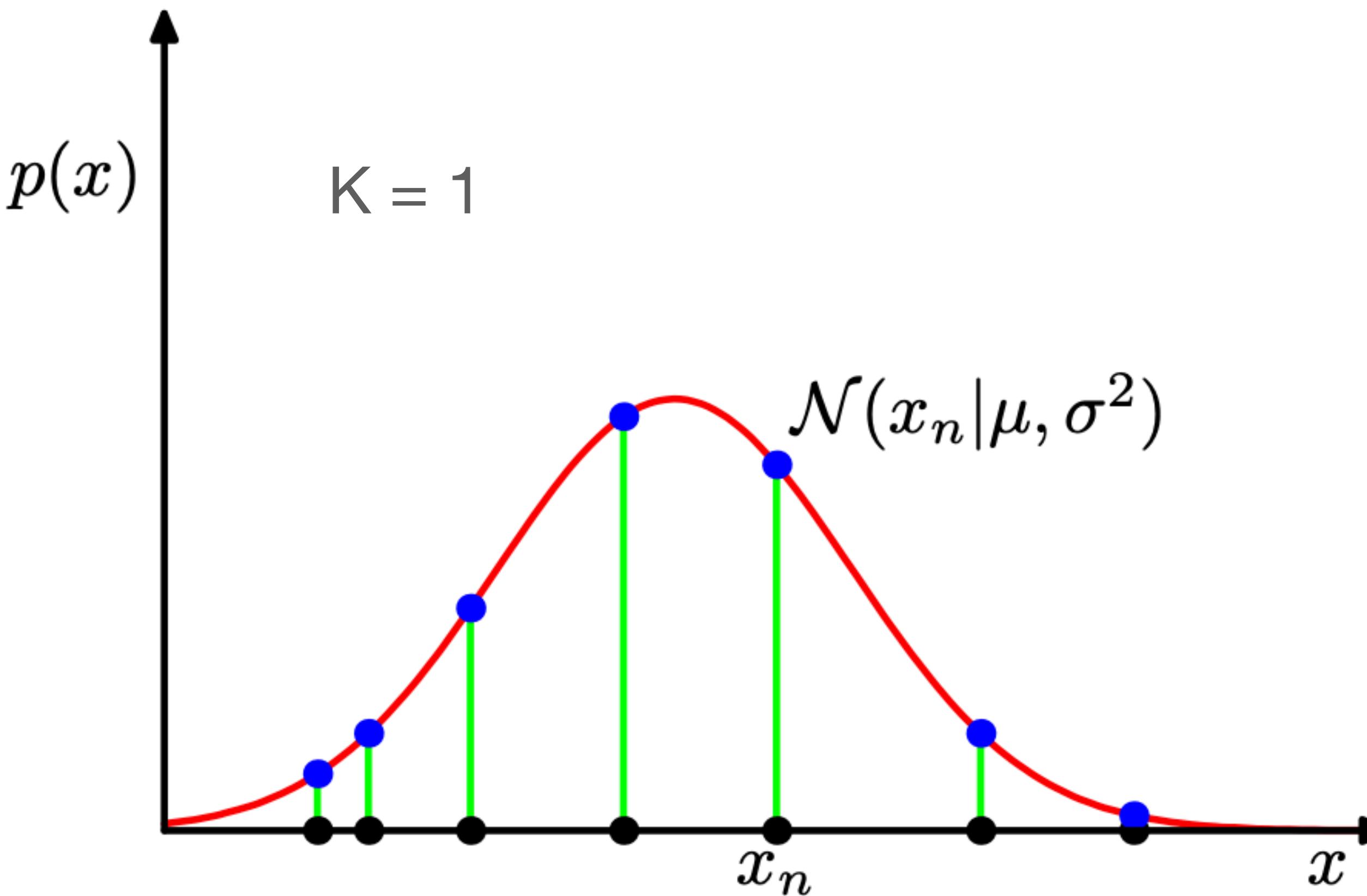
3. K-means as a special case

4. Theoretical foundation\*: EM algorithm (1 is E step, 2 is M step)

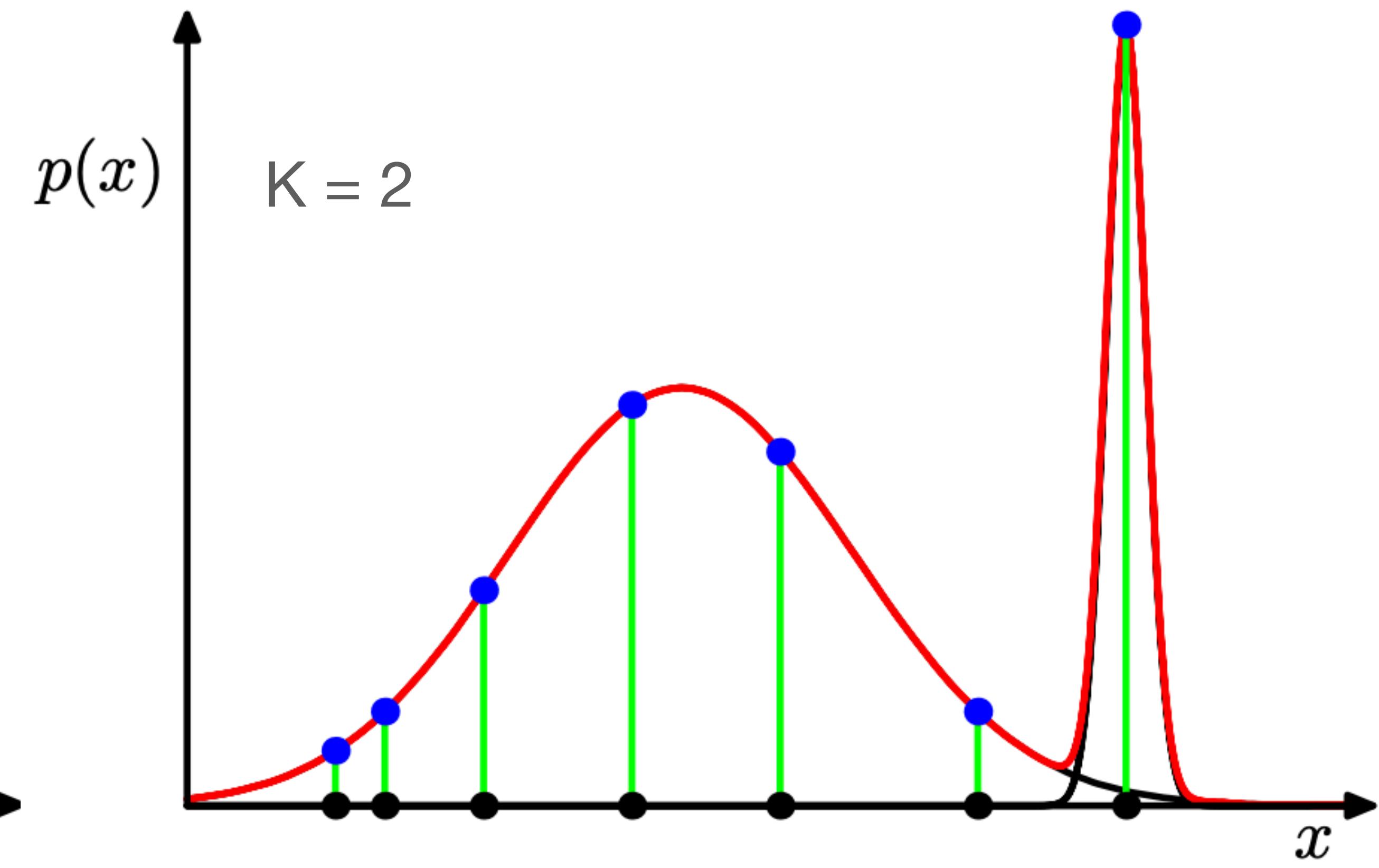
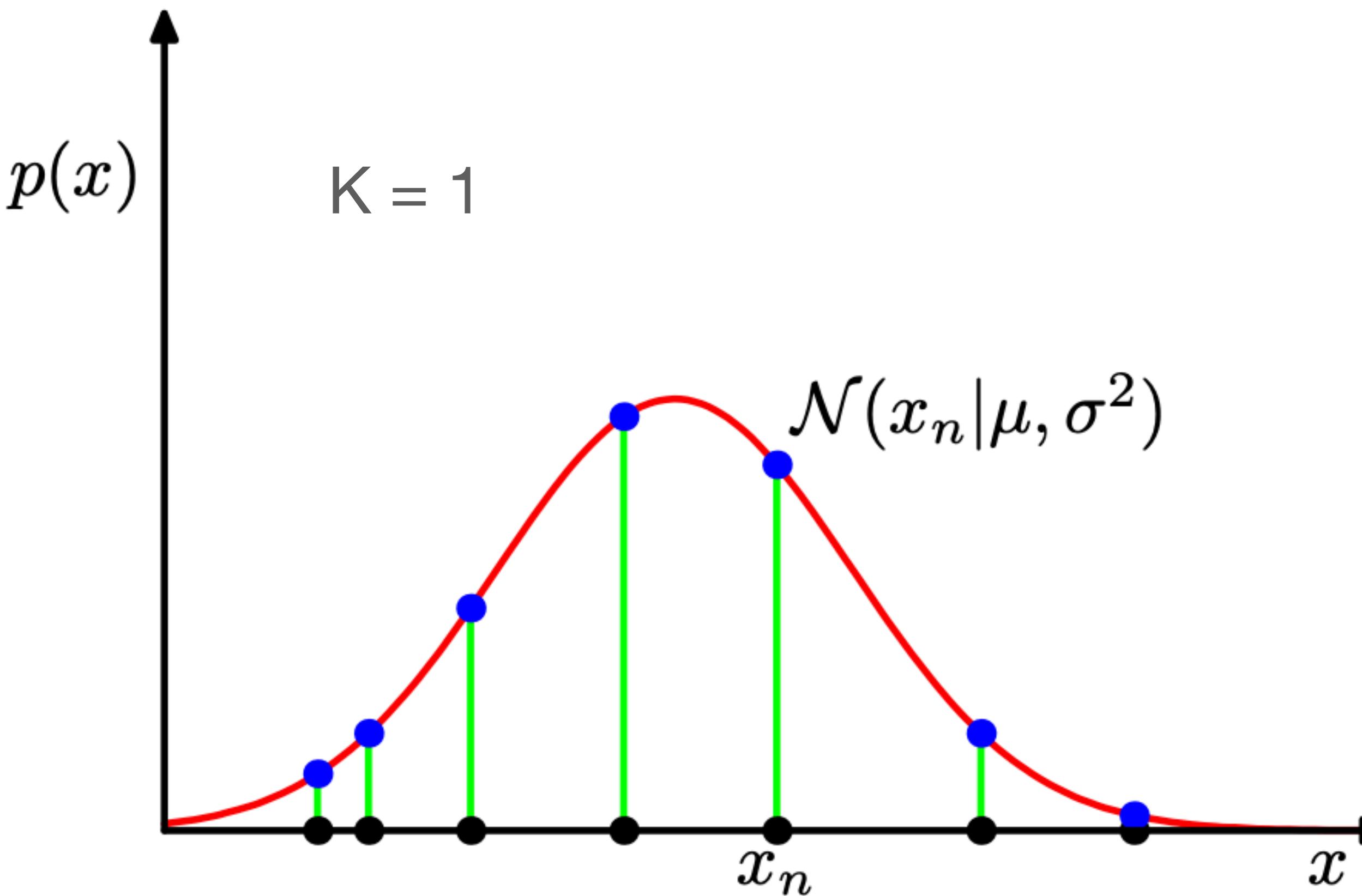
# Maximum likelihood - overfitting



# Maximum likelihood - overfitting



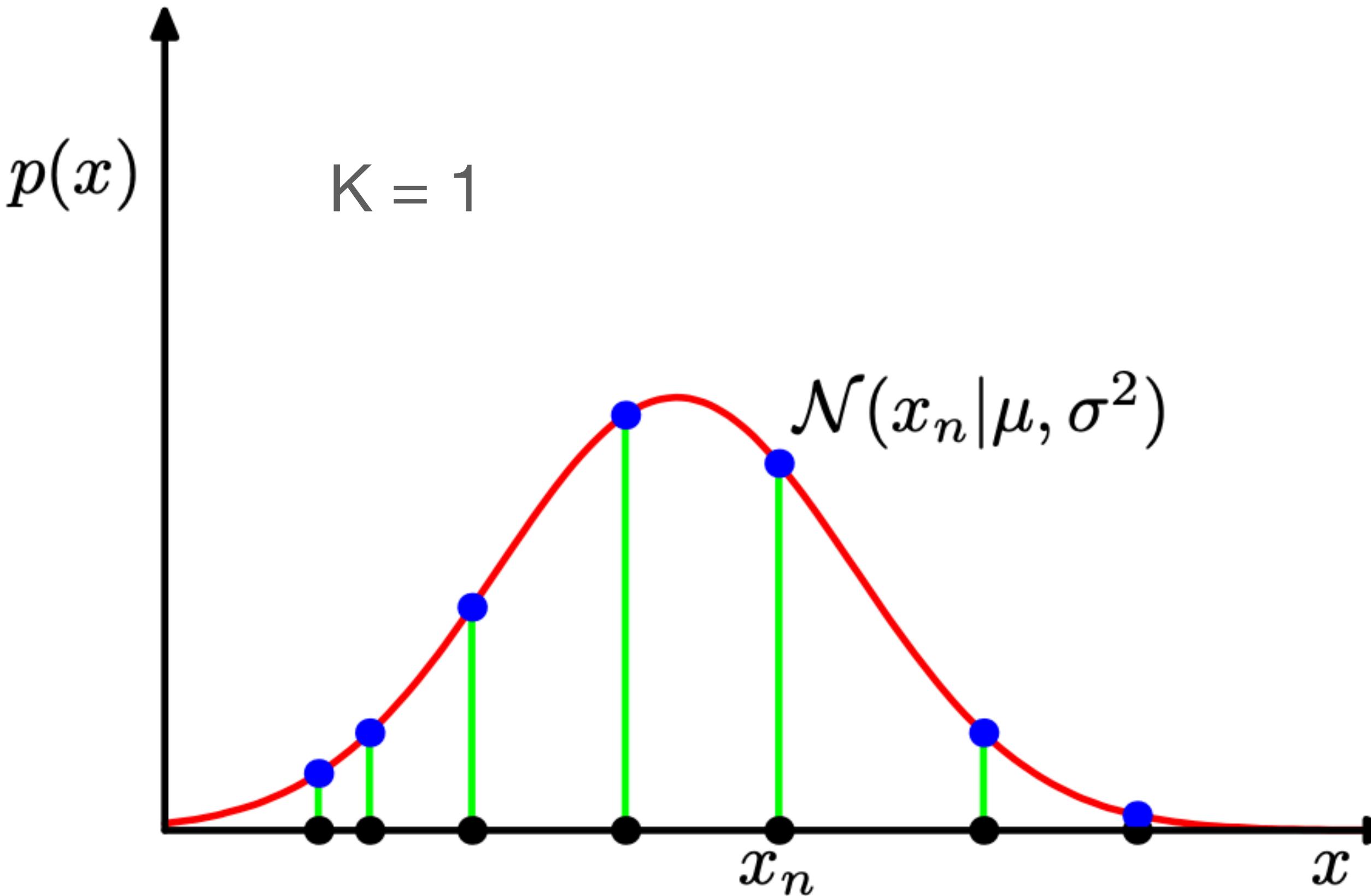
# Maximum likelihood - overfitting



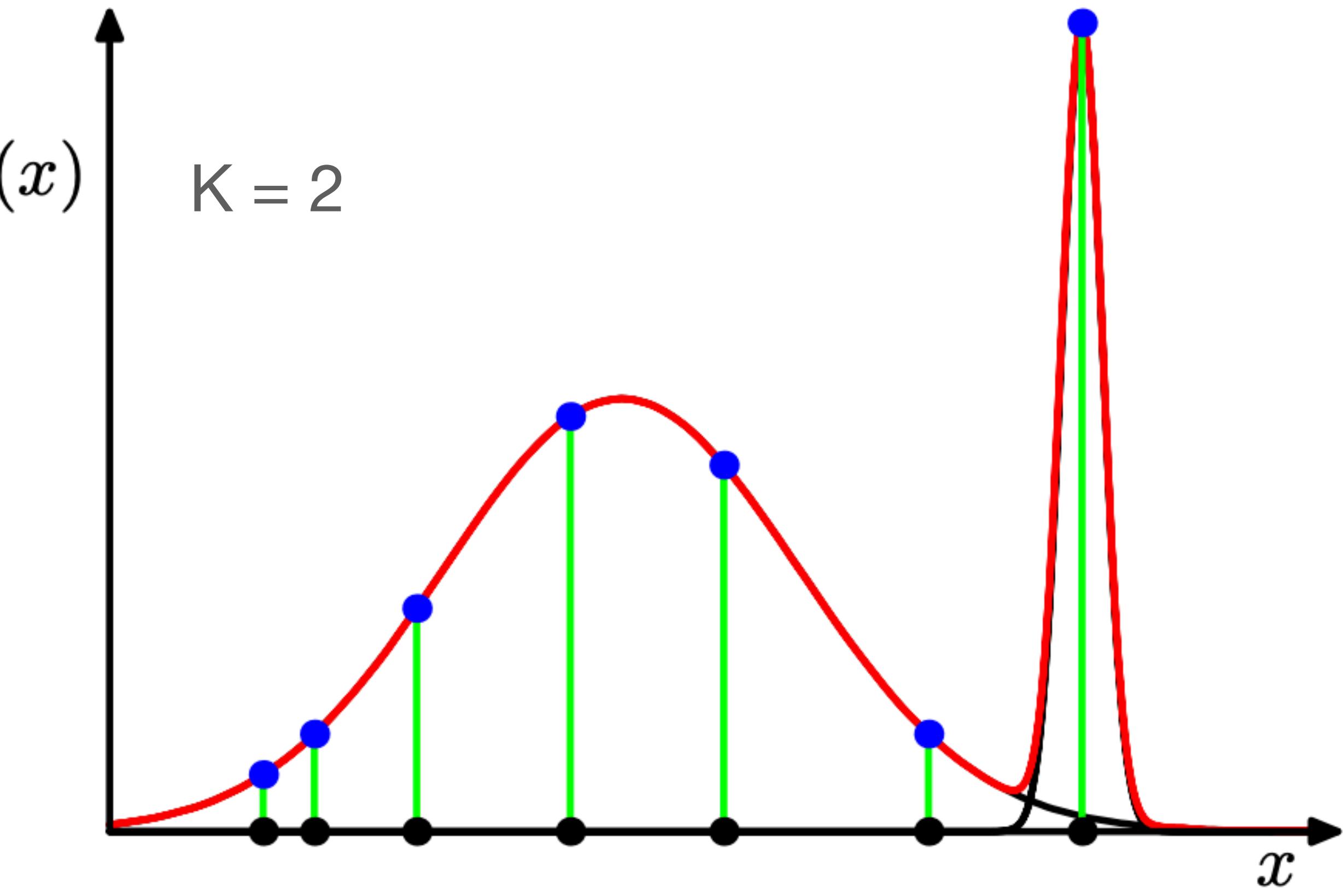
$$p_2(x_7 | \mu_2, \sigma_2^2) = \mathcal{N}(x_7; x_7, \sigma_2^2) = \frac{1}{\sqrt{2\pi}} \frac{1}{\sigma_2}$$

One component collapses to a data point!

# Maximum likelihood - overfitting



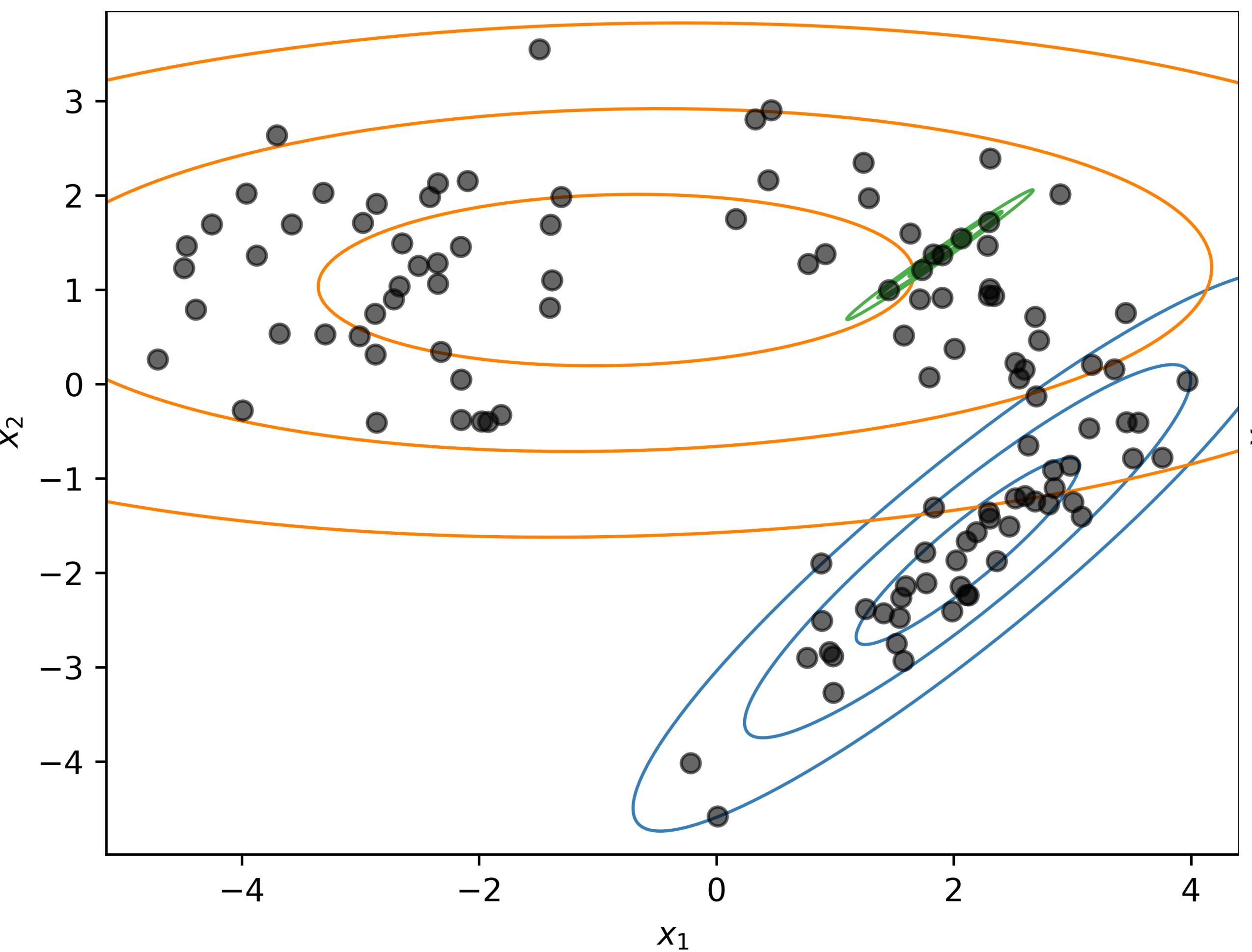
**Hot-fix:** reset the mean to a randomly chosen value, and the covariance to some large value and hope we don't have to reset again \\_(\_)



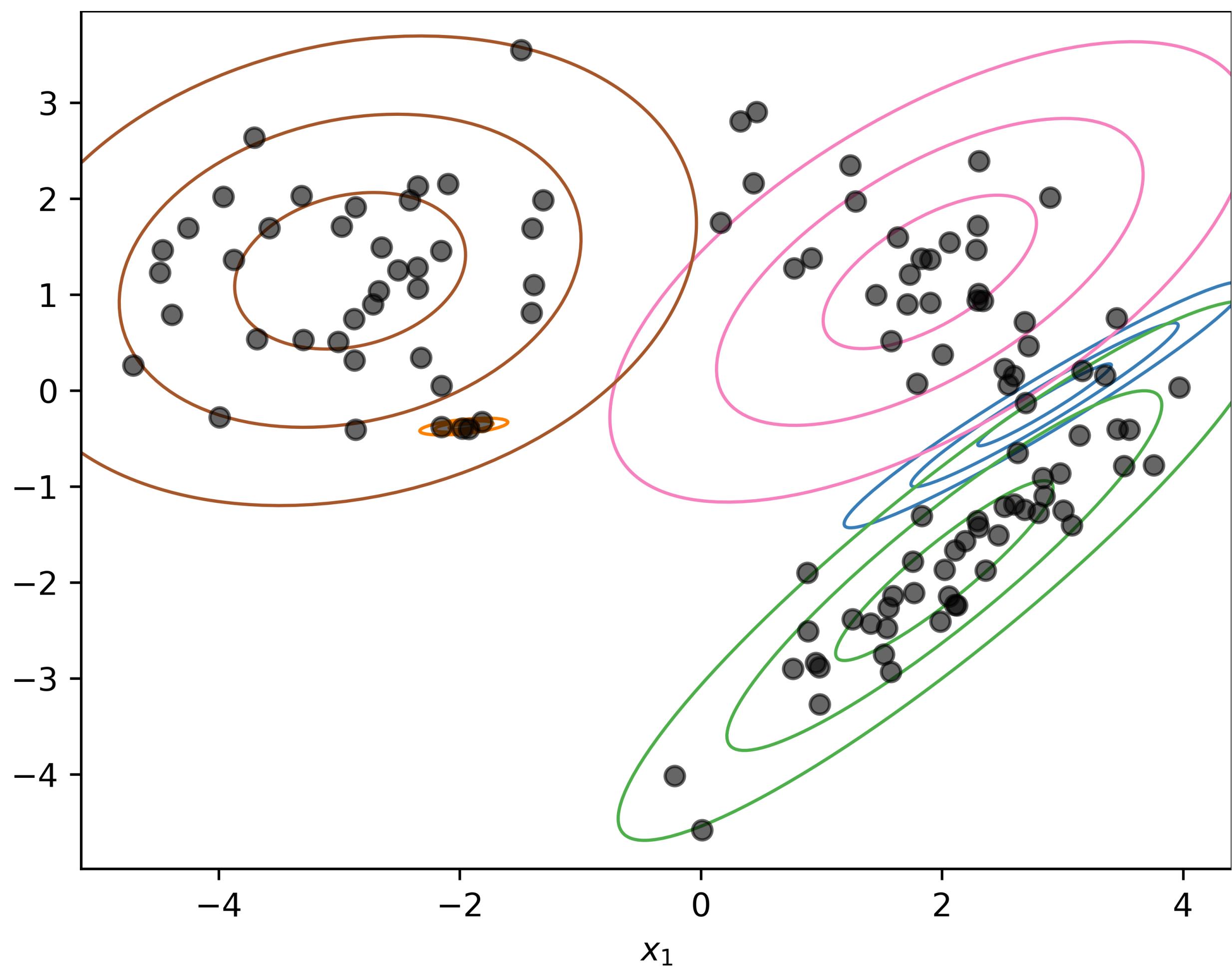
$$p_2(x_7 | \mu_2, \sigma_2^2) = \mathcal{N}(x_7; x_7, \sigma_2^2) = \frac{1}{\sqrt{2\pi}} \frac{1}{\sigma_2}$$

One component collapses to a data point!

# Local minima

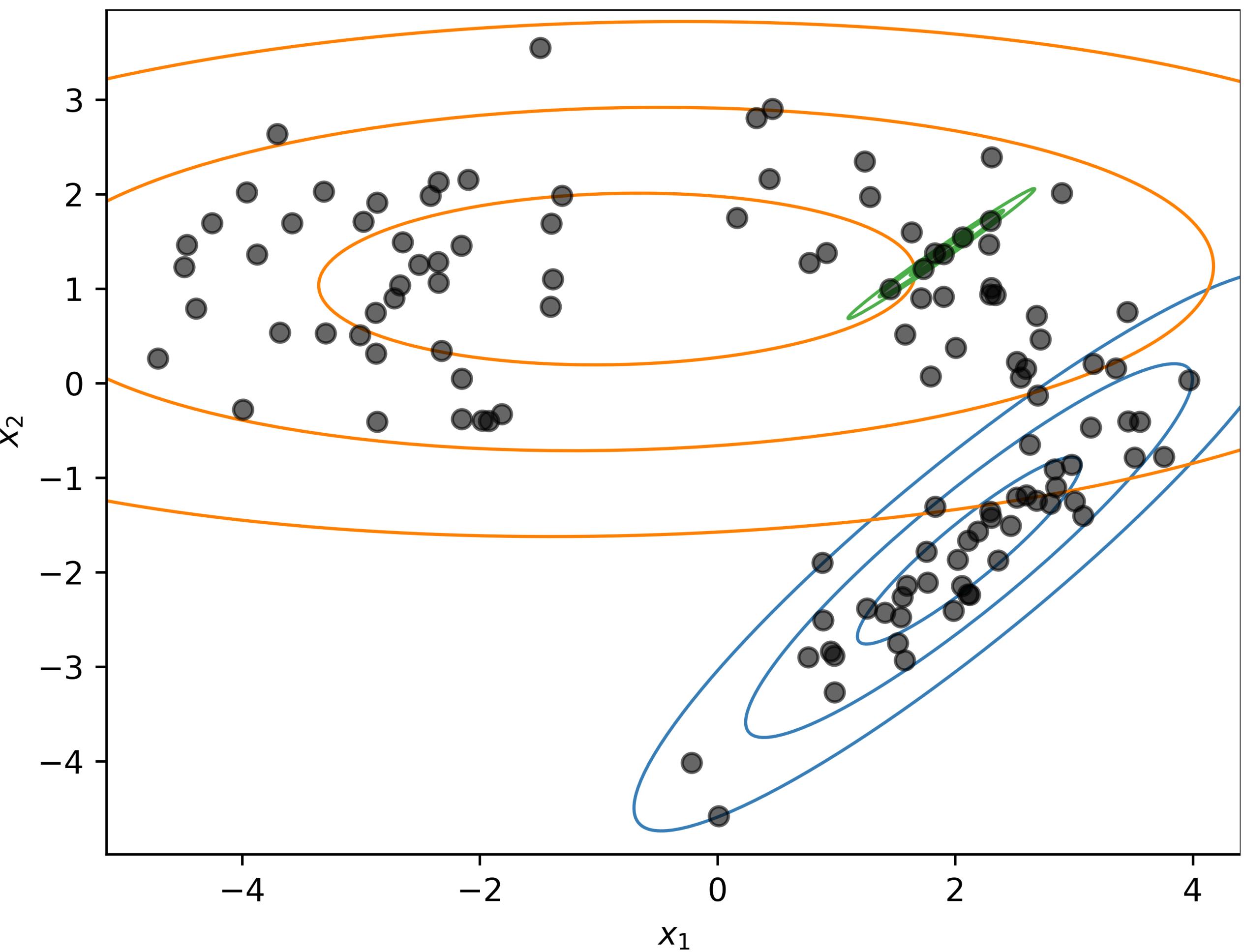


$K = 3$

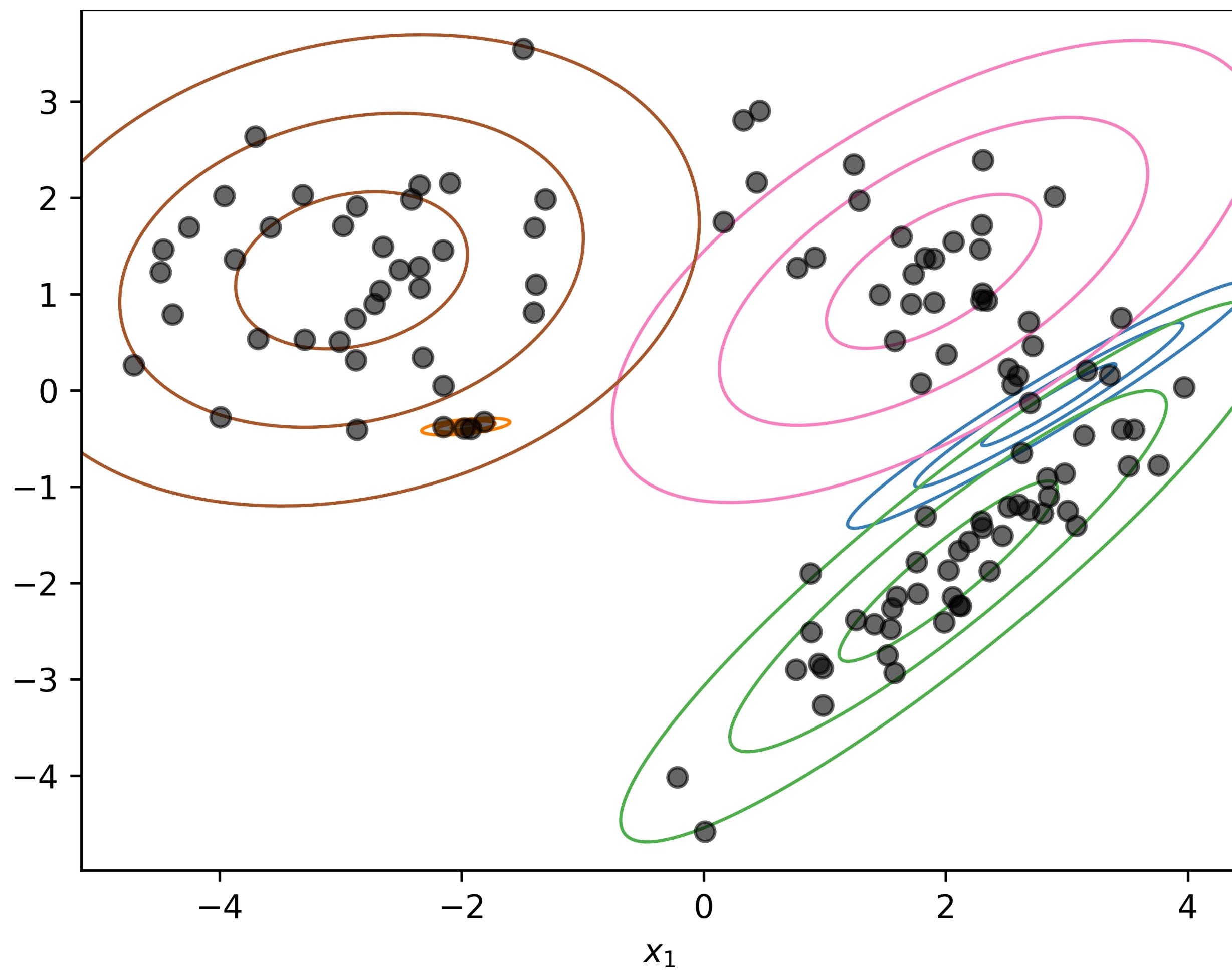


$K = 5$

# Local minima



$K = 3$



$K = 5$

Practical initialisation: k-means, mean/cov = sample mean/cov, mixture weights = proportions

# Overview

Last lecture:

1. **Model:** Gaussian mixture models
2. **Learning:** parameter learning using *maximum likelihood*
3. **Iterative** procedure, derive mean update

**Initialise**  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

**Repeat** until convergence:

1. **Evaluate responsibilities**  $r_{nk}$  for every data point  $x_n$  using current parameters:

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}$$

2. **Re-estimate parameters**  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$  using the current responsibilities  $r_{nk}$ :

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n, \quad N_k = \sum_{n=1}^N r_{nk}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^\top$$

$$\pi_k = \frac{N_k}{N}$$

This lecture:

1. Derive covariance and mixture weight updates
2. Issues with maximum likelihood and workarounds
3. K-means as a special case
4. Theoretical foundation\*: EM algorithm (1 is E step, 2 is M step)

# Relation to k-means

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}$$

Assume  $\Sigma_1 = \Sigma_2 = \dots = \Sigma_K = \epsilon \mathbf{I}$

Numerical example:  $x_n = 5$ ,  $K = 3$

# Relation to k-means

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}$$

Assume  $\Sigma_1 = \Sigma_2 = \dots = \Sigma_K = \epsilon \mathbf{I}$

Numerical example:  $x_n = 5$ ,  $K = 3$

$\mu_k$	$\pi_k$	$\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$	$r_{nk}$
8	0.3	0.02	0.27
-2	0.4	0.00016	0.01
4	0.3	0.0528	0.72

Soft assignment

$$\epsilon = 4$$

# Relation to k-means

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}$$

Assume  $\Sigma_1 = \Sigma_2 = \dots = \Sigma_K = \epsilon \mathbf{I}$

Numerical example:  $x_n = 5$ ,  $K = 3$

$\mu_k$	$\pi_k$	$\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$	$r_{nk}$
8	0.3	0.02	0.27
-2	0.4	0.00016	0.01
4	0.3	0.0528	0.72

Soft assignment

$$\epsilon = 4$$

$\mu_k$	$\pi_k$	$\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$	$r_{nk}$
8	0.3	$0.44 \times 10^{-195}$	0
-2	0.4	$1.5 \times 10^{-1064}$	0
4	0.3	$0.23 \times 10^{-22}$	1

Hard assignment

$$\epsilon = 0.01$$

# Overview

Last lecture:

1. **Model:** Gaussian mixture models
2. **Learning:** parameter learning using *maximum likelihood*
3. **Iterative** procedure, derive mean update

**Initialise**  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

**Repeat** until convergence:

1. **Evaluate responsibilities**  $r_{nk}$  for every data point  $x_n$  using current parameters:

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}$$

2. **Re-estimate parameters**  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$  using the current responsibilities  $r_{nk}$ :

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n, \quad N_k = \sum_{n=1}^N r_{nk}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^\top$$

$$\pi_k = \frac{N_k}{N}$$

This lecture:

1. Derive covariance and mixture weight updates
2. Issues with maximum likelihood and workarounds
3. K-means as a special case
4. Theoretical foundation\*: EM algorithm (1 is E step, 2 is M step)

# GMMs - iterative procedure = *EM* algorithm

Initialise  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

Repeat until convergence:

1. Evaluate responsibilities  $r_{nk}$  for every data point  $x_n$  using current parameters:

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}$$

**E step**

2. Re-estimate parameters  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$  using the current responsibilities  $r_{nk}$ :

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n, \quad N_k = \sum_{n=1}^N r_{nk}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^\top$$

$$\pi_k = \frac{N_k}{N}$$

**M step**

# Expectation maximisation (EM) algorithm\*

$$\text{Log-likelihood: } L(\theta) = \log p(X | \theta) = \sum_{n=1}^N \log p(x_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

**Goal:** maximum likelihood,  $\operatorname{argmax}_\theta p(X | \theta)$ , which is equivalent to  $\operatorname{argmax}_\theta \log p(X | \theta)$

Assume: **latent/unknown/hidden variables**  $\{z_n\}_{n=1}^N$  representing *component assignment*

$$L(\theta) = \log p(X | \theta) = \log \int p(X, z | \theta) dz = \log \int q(z) \frac{p(X, z | \theta)}{q(z)} dz \geq \int q(z) \log \frac{p(X, z | \theta)}{q(z)} dz := \mathcal{F}(q(z), \theta)$$

An arbitrary  $q(z)$       Jensen's inequality      Lower bound on  $L(\theta)$

Instead of maximising  $L(\theta)$  directly, we will maximise the lower bound  $\mathcal{F}(q(z), \theta)$ , alternating between  $q(z)$  and  $\theta$  while keeping the other fixed.

# EM for GMMs

1. Evaluate responsibilities  $r_{nk}$  for every data point  $x_n$  using current parameters:

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}$$

**E step**

2. Re-estimate parameters  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$  using the current responsibilities  $r_{nk}$ :

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n, \quad N_k = \sum_{n=1}^N r_{nk}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^\top$$

$$\pi_k = \frac{N_k}{N}$$

**M step**

# EM for GMMs

Maximising lower bound wrt  $q(z)$   
whilst keeping  $\theta$  fixed

$$q(z)^{(t+1)} \leftarrow \operatorname{argmax}_{q(z)} \mathcal{F}(q(z), \theta^{(t)}) = p(z | X, \theta^{(t)})$$

fill in values for the hidden variables according to their posterior probabilities

1. **Evaluate responsibilities**  $r_{nk}$  for every data point  $x_n$  using current parameters:

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}$$

**E step**

2. **Re-estimate parameters**  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$  using the current responsibilities  $r_{nk}$ :

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n, \quad N_k = \sum_{n=1}^N r_{nk}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^\top$$

$$\pi_k = \frac{N_k}{N}$$

**M step**

# EM for GMMs

Maximising lower bound wrt  $q(z)$   
whilst keeping  $\theta$  fixed

$$q(z)^{(t+1)} \leftarrow \operatorname{argmax}_{q(z)} \mathcal{F}(q(z), \theta^{(t)}) = p(z | X, \theta^{(t)})$$

fill in values for the hidden variables according to their posterior probabilities

1. **Evaluate responsibilities**  $r_{nk}$  for every data point  $x_n$  using current parameters:

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}$$

**E step**

2. **Re-estimate parameters**  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$  using the current responsibilities  $r_{nk}$ :

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n, \quad N_k = \sum_{n=1}^N r_{nk}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^\top$$

$$\pi_k = \frac{N_k}{N}$$

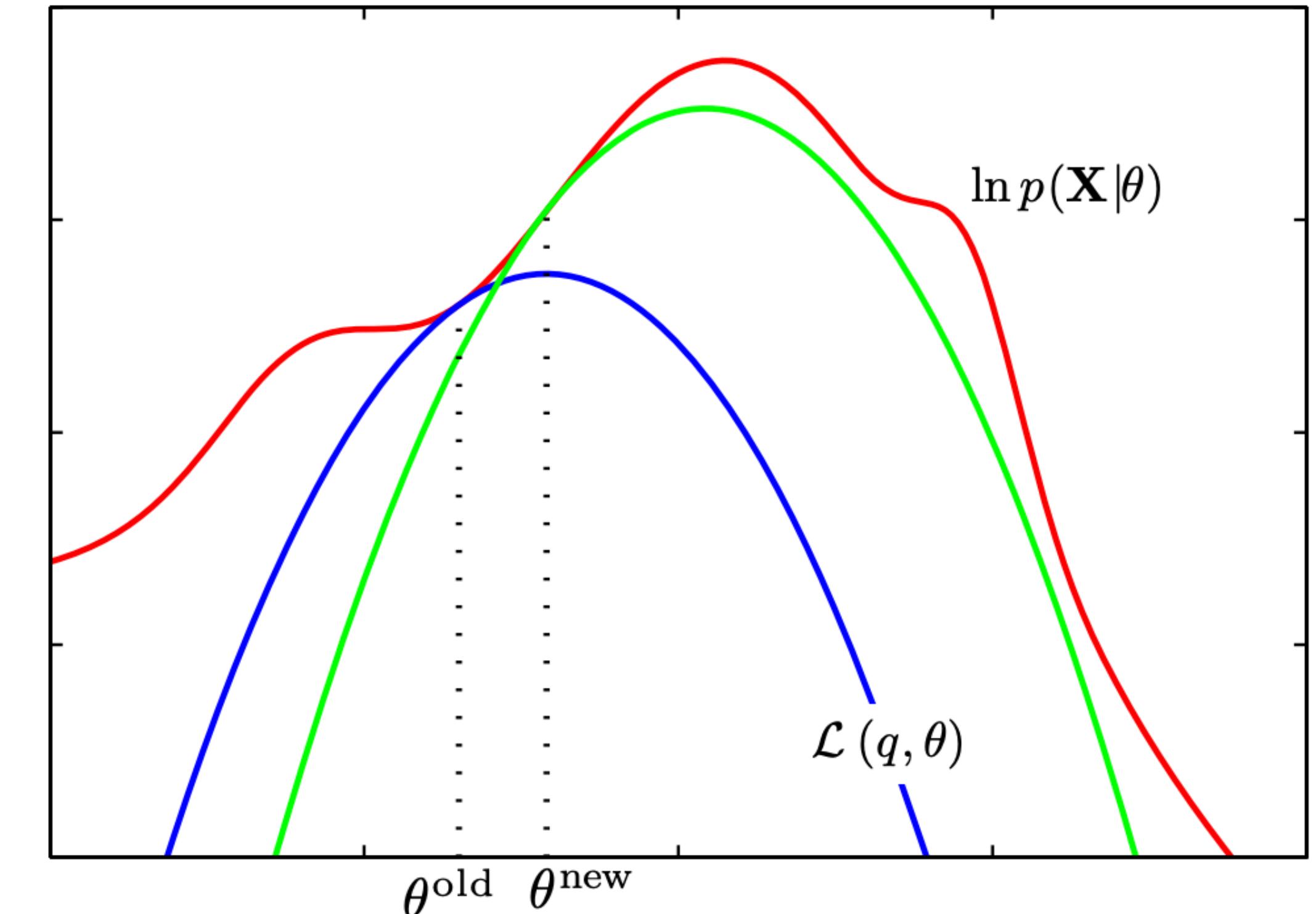
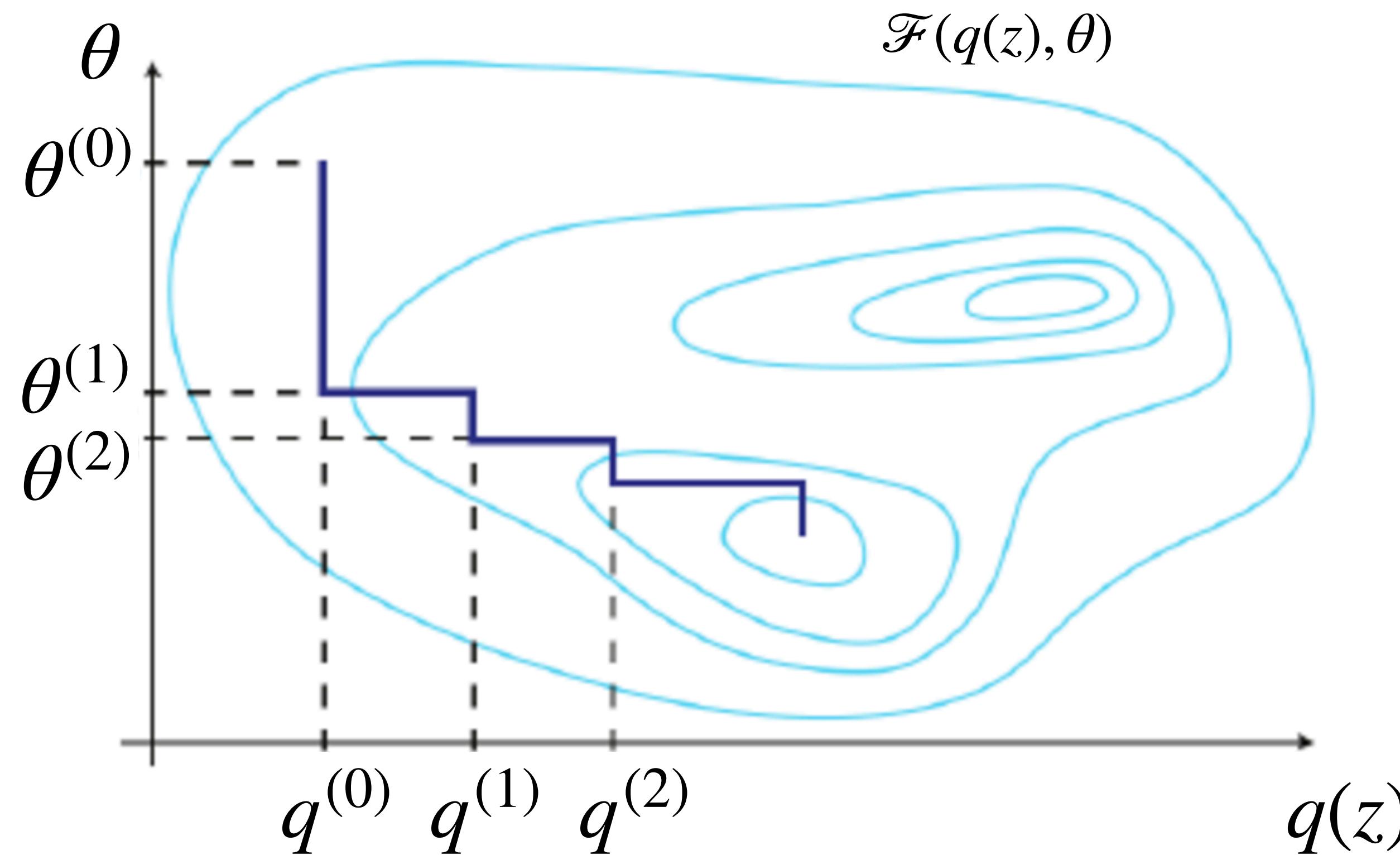
**M step**

Maximising lower bound wrt  $\theta$  whilst  
keeping  $q(z)$  fixed

$$\theta^{(t+1)} \leftarrow \operatorname{argmax}_\theta \mathcal{F}(q(z)^{t+1}, \theta)$$

learn model as if hidden variables were not hidden

# EM algorithm - visualisation



# EM algorithm - more visualisation

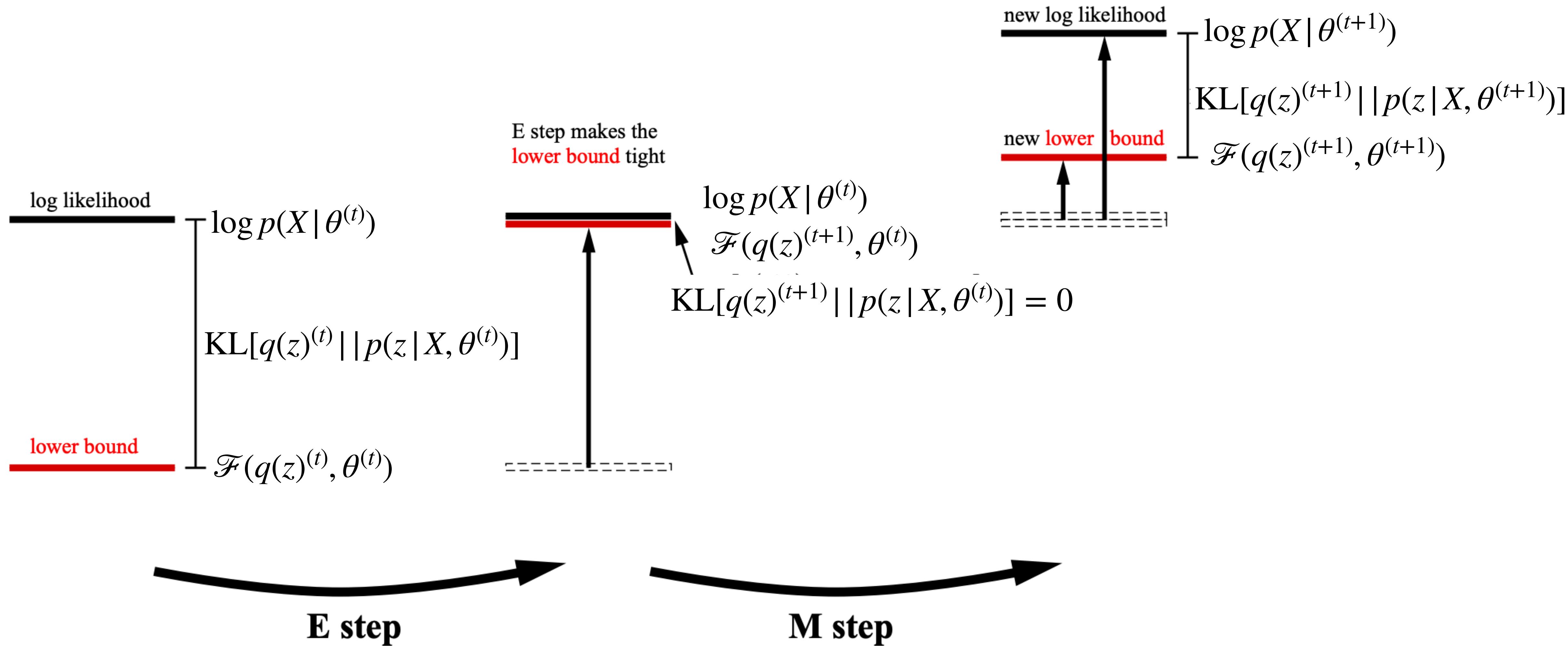


Figure credit: Beal (2003), Variational Algorithms for Approximate Bayesian Inference

# Gaussian mixture models: check your understanding

**Model:** 
$$p_{\theta}(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

where we defined  $\theta = \{\pi_k, \mu_k, \Sigma_k \text{ for } k = 1, 2, \dots, K\}$  as the parameters of the GMM.

- If  $K$  takes a greater value, the likelihood becomes greater after convergence.
- Assume we have  $N$  data points. The maximum likelihood will be achieved if we set  $K = N$ .
- In GMM, the EM algorithm gives us *global* minimum, because we can update  $\pi_k$ ,  $\mu_k$ , and  $\Sigma_k$  through closed-form solutions.
- GMM has a higher computational complexity than k-means.

# Gaussian mixture models: check your understanding

**Model:** 
$$p_{\theta}(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

where we defined  $\theta = \{\pi_k, \mu_k, \Sigma_k \text{ for } k = 1, 2, \dots, K\}$  as the parameters of the GMM.

- If  $K$  takes a greater value, the likelihood becomes greater after convergence. Yes
- Assume we have  $N$  data points. The maximum likelihood will be achieved if we set  $K = N$ .
- In GMM, the EM algorithm gives us *global* minimum, because we can update  $\pi_k$ ,  $\mu_k$ , and  $\Sigma_k$  through closed-form solutions.
- GMM has a higher computational complexity than k-means.

# Gaussian mixture models: check your understanding

**Model:** 
$$p_{\theta}(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

where we defined  $\theta = \{\pi_k, \mu_k, \Sigma_k \text{ for } k = 1, 2, \dots, K\}$  as the parameters of the GMM.

- If  $K$  takes a greater value, the likelihood becomes greater after convergence. Yes
- Assume we have  $N$  data points. The maximum likelihood will be achieved if we set  $K = N$ . Yes
- In GMM, the EM algorithm gives us *global* minimum, because we can update  $\pi_k$ ,  $\mu_k$ , and  $\Sigma_k$  through closed-form solutions.
- GMM has a higher computational complexity than k-means.

# Gaussian mixture models: check your understanding

**Model:** 
$$p_{\theta}(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

where we defined  $\theta = \{\pi_k, \mu_k, \Sigma_k \text{ for } k = 1, 2, \dots, K\}$  as the parameters of the GMM.

- If  $K$  takes a greater value, the likelihood becomes greater after convergence. Yes
- Assume we have  $N$  data points. The maximum likelihood will be achieved if we set  $K = N$ . Yes
- In GMM, the EM algorithm gives us *global* minimum, because we can update  $\pi_k$ ,  $\mu_k$ , and  $\Sigma_k$  through closed-form solutions. No
- GMM has a higher computational complexity than k-means.

# Gaussian mixture models: check your understanding

**Model:** 
$$p_{\theta}(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \quad 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1$$

where we defined  $\theta = \{\pi_k, \mu_k, \Sigma_k \text{ for } k = 1, 2, \dots, K\}$  as the parameters of the GMM.

- If  $K$  takes a greater value, the likelihood becomes greater after convergence. Yes
- Assume we have  $N$  data points. The maximum likelihood will be achieved if we set  $K = N$ . Yes
- In GMM, the EM algorithm gives us *global* minimum, because we can update  $\pi_k$ ,  $\mu_k$ , and  $\Sigma_k$  through closed-form solutions. No
- GMM has a higher computational complexity than k-means. Yes

# Overview

Last lecture:

1. **Model:** Gaussian mixture models
2. **Learning:** parameter learning using *maximum likelihood*
3. **Iterative** procedure, derive mean update

**Initialise**  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

**Repeat** until convergence:

1. **Evaluate responsibilities**  $r_{nk}$  for every data point  $x_n$  using current parameters:

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n; \mu_j, \Sigma_j)}$$

2. **Re-estimate parameters**  $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$  using the current responsibilities  $r_{nk}$ :

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n, \quad N_k = \sum_{n=1}^N r_{nk}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^\top$$

$$\pi_k = \frac{N_k}{N}$$

This lecture:

1. Derive covariance and mixture weight updates
2. Issues with maximum likelihood and workarounds
3. K-means as a special case
4. Theoretical foundation\*: EM algorithm (1 is E step, 2 is M step)

\* not in assignment/exam