


# Probabilistic Modelling Principles

**Yingzhen Li**

Department of Computing  
Imperial College London

@liyzhen2  
yingzhen.li@imperial.ac.uk

October 21, 2022

# Principles of probabilistic modelling

Have you ever wondered about the following questions:

- Why using  $\ell_2$  loss in many regression problems?
- Where does the cross-entropy loss come from?
- What is a good principle for choosing a good loss function?

# Principles of probabilistic modelling

Have you ever wondered about the following questions:

- Why using  $\ell_2$  loss in many regression problems?
- Where does the cross-entropy loss come from?
- What is a good principle for choosing a good loss function?

**Probabilistic modelling** gives you good answers for all of them!

# Principles of probabilistic modelling

Probabilistic modelling is about:

1. making model assumptions on **how the data is generated**
2. estimating model parameters under probabilistic principles
3. model checking using data, and repeat 1 - 3
4. using the fitted model for downstream tasks

# Example: coin flipping

Imagine you'd like to predict the next coin flip result:



- Assume  $x_1, x_2, \dots, x_N$  are observed **independent** coin flip results using the **same** coin,
- I.e.,  $x_1, \dots, x_N$  are sampled i.i.d. from the same **data distribution**  $\pi(x)$
- However, we don't know  $\pi(x)$

# Example: coin flipping

Imagine you'd like to predict the next coin flip result:



Probabilistic modelling is about:

1. Assume  $x$  is sampled from  $p(x|\theta)$   $\Leftarrow$  **our probabilistic model**
2. estimating  $\theta$  under probabilistic principles  
such as MLE, MAP, posterior inference  $\Leftarrow$  **learning the model**
3. check if  $p(x|\theta^*)$  fits  $\pi(x)$  well, and repeat 1 - 3  $\Leftarrow$  **model checking**
4. making prediction for next coin flip result using  $p(x|\theta^*)$

# Example: coin flipping

Imagine you'd like to predict the next coin flip result:



Step 1: Assume  $x$  is sampled from  $p(x|\theta)$

$$x = \begin{cases} 1, & \text{with probability } \theta \\ 0, & \text{with probability } 1 - \theta \end{cases}, \quad \theta \in [0, 1].$$

$$\Leftrightarrow p(x|\theta) = \text{Bern}(\theta).$$

- Likelihood of  $\theta$  given observed  $x$ :  $\ell(\theta) = p(x|\theta)$

# Example: coin flipping

Imagine you'd like to predict the next coin flip result:



Step 2: estimating  $\theta$  using probabilistic principles

Here we consider **maximum likelihood estimation (MLE)**

Idea of MLE: for datapoints  $x$  sampled from  $\pi(x)$

- We want to find  $\theta^*$  such that  $p(x|\theta^*) \approx \pi(x)$



## Example: coin flipping

Imagine you'd like to predict the next coin flip result:



Step 2: estimating  $\theta$  using probabilistic principles

Here we consider **maximum likelihood estimation (MLE)**

Idea of MLE: for datapoints  $x$  sampled from  $\pi(x)$

- ▶ We want to find  $\theta^*$  such that  $p(x|\theta^*) \approx \pi(x)$
- ▶ We need to measure the “closeness” of the two distributions  $\Rightarrow$  use the KL divergence

$$\text{KL}[\pi(x)||p(x|\theta)] = \mathbb{E}_{\pi(x)} \left[ \log \frac{\pi(x)}{p(x|\theta)} \right]$$

# Example: coin flipping

Imagine you'd like to predict the next coin flip result:



Step 2: estimating  $\theta$  using probabilistic principles

Here we consider **maximum likelihood estimation (MLE)**

Idea of MLE: for datapoints  $x$  sampled from  $\pi(x)$

- ▶ We want to find  $\theta^*$  such that  $p(x|\theta^*) \approx \pi(x)$
- ▶ We want this KL to be small:

$$\theta^* = \arg \min_{\theta} \text{KL}[\pi(x) || p(x|\theta)]$$

## Example: coin flipping

Imagine you'd like to predict the next coin flip result:



Step 2: estimating  $\theta$  using probabilistic principles

Here we consider **maximum likelihood estimation (MLE)**

Idea of MLE: for datapoints  $x$  sampled from  $\pi(x)$

- We want to find  $\theta^*$  such that  $p(x|\theta^*) \approx \pi(x)$

$$\Leftrightarrow \theta^* = \arg \max_{\theta} \mathbb{E}_{\pi(x)} [\log p(x|\theta)]$$

## Example: coin flipping

Imagine you'd like to predict the next coin flip result:



Step 2: estimating  $\theta$  using probabilistic principles

Here we consider **maximum likelihood estimation (MLE)**

Idea of MLE: for datapoints  $x$  sampled from  $\pi(x)$

- ▶ We want to find  $\theta^*$  such that  $p(x|\theta^*) \approx \pi(x)$
- ▶ Estimate using dataset  $\mathcal{D} = \{x_1, \dots, x_N\}$  sampled from  $\pi(x)$ :

$$\theta^* = \arg \max_{\theta} \frac{1}{N} \sum_{n=1}^N \log p(x_n|\theta)$$

# Example: coin flipping

Imagine you'd like to predict the next coin flip result:



Step 2: estimating  $\theta$  using probabilistic principles

Here we consider **maximum likelihood estimation (MLE)**

- Estimate using dataset  $\mathcal{D} = \{x_1, \dots, x_N\}$  sampled from  $\pi(x)$ :

$$\theta^* = \arg \max_{\theta} \frac{1}{N} \sum_{n=1}^N \log p(x_n | \theta)$$

- model assumption:  $p(x | \theta) = \text{Bern}(\theta)$

$$\Rightarrow \theta^* = \arg \max_{\theta} \frac{1}{N} \sum_{n=1}^N x_n \log \theta + (1 - x_n) \log(1 - \theta)$$

## Example: coin flipping

Imagine you'd like to predict the next coin flip result:



Step 2: estimating  $\theta$  using probabilistic principles

Here we consider **maximum likelihood estimation (MLE)**

- Estimate using dataset  $\mathcal{D} = \{x_1, \dots, x_N\}$  sampled from  $\pi(x)$ :

$$\theta^* = \arg \max_{\theta} \frac{1}{N} \sum_{n=1}^N \log p(x_n | \theta)$$

- solution by zeroing the gradient:

$$\frac{1}{N} \sum_{n=1}^N x_n \theta^{-1} - (1 - x_n)(1 - \theta)^{-1} = 0 \quad \Rightarrow \quad \theta^* = \frac{1}{N} \sum_{n=1}^N x_n$$

## Example: coin flipping

Imagine you'd like to predict the next coin flip result:



Step 3: check if  $p(x|\theta^*)$  fits  $\pi(x)$  well  
(We assume the model has passed here)

## Example: coin flipping

Imagine you'd like to predict the next coin flip result:



Step 4: making prediction for next coin flip result using  $p(x|\theta^*)$

$$\theta^* = \frac{1}{N} \sum_{x \in \mathcal{D}} x$$

$$\Rightarrow x_{N+1} = \begin{cases} 1, & \text{with probability } \frac{1}{N} \sum_{n=1}^N x_n \\ 0, & \text{with probability } 1 - \frac{1}{N} \sum_{n=1}^N x_n \end{cases} .$$



# Principles of probabilistic modelling

Datapoints  $(x, y)$  are sampled from an **unknown ground truth distribution**  $\pi(x, y)$

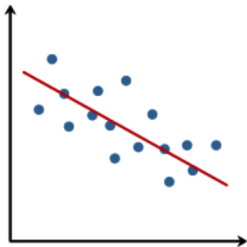
Probabilistic modelling is about (in supervised learning case):

1. Assuming the output  $y$  given  $x$  is sampled from

$$p(y|x, \theta)$$

2. estimating  $\theta$  under probabilistic principles such as MLE, MAP, posterior inference
3. check if  $p(y|x, \theta^*)$  fits  $\pi(y|x)$  well, and repeat 1 - 3
4. using  $p(y|x, \theta^*)$  for predictions

# Probabilistic modelling: linear regression

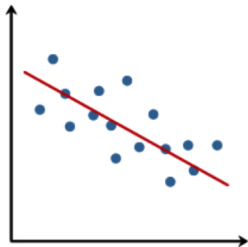


Linear regression

$$f(x, \theta) = \theta^\top x,$$

$$y = f(x, \theta) + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$$

# Probabilistic modelling: linear regression

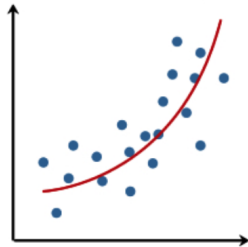


Linear regression

$$f(x, \theta) = \theta^\top x,$$

$$y = f(x, \theta) + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$\Rightarrow$



Non-linear regression

$$f(x, \theta) = \theta^\top \phi(x)$$

$$y = f(x, \theta) + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$$

# Probabilistic modelling: linear regression

Step 1: making assumptions about the output generation process

$$y = \boldsymbol{\theta}^\top \boldsymbol{\phi}(x) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

- $\boldsymbol{\theta}$  is the model parameter
- $\boldsymbol{\phi}(x)$  is a pre-defined feature mapping (e.g., polynomial features)

# Probabilistic modelling: linear regression

Step 1: making assumptions about the output generation process

$$y = \boldsymbol{\theta}^\top \boldsymbol{\phi}(x) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

- $\boldsymbol{\theta}$  is the model parameter
- $\boldsymbol{\phi}(x)$  is a pre-defined feature mapping (e.g., polynomial features)

Probabilistic formulation:

- The distribution of  $y$  given  $x$  under model assumption:

$$p(y|x, \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}^\top \boldsymbol{\phi}(x), \sigma^2)$$

- Likelihood of  $\boldsymbol{\theta}$  given observed data  $(x, y)$ :

$$\ell(\boldsymbol{\theta}) = p(y|x, \boldsymbol{\theta})$$

# Probabilistic modelling: linear regression

Step 2: estimating  $\theta$  using **maximum likelihood estimation (MLE)**

Idea of MLE: for datapoints  $(x, y)$  sampled from  $\pi(x, y)$

- We want to find  $\theta^*$  such that  $p(y|x, \theta^*) \approx \pi(y|x)$

# Probabilistic modelling: linear regression

Step 2: estimating  $\theta$  using **maximum likelihood estimation (MLE)**

Idea of MLE: for datapoints  $(x, y)$  sampled from  $\pi(x, y)$

- ▶ We want to find  $\theta^*$  such that  $p(y|x, \theta^*) \approx \pi(y|x)$
- ▶ We need to measure the “closeness” of the two distributions  $\Rightarrow$  use the KL divergence

$$\text{KL}[\pi(y|x) || p(y|x, \theta)] = \mathbb{E}_{\pi(y|x)} \left[ \log \frac{\pi(y|x)}{p(y|x, \theta)} \right]$$

# Probabilistic modelling: linear regression

Step 2: estimating  $\theta$  using **maximum likelihood estimation (MLE)**

Idea of MLE: for datapoints  $(x, y)$  sampled from  $\pi(x, y)$

- ▶ We want to find  $\theta^*$  such that  $p(y|x, \theta^*) \approx \pi(y|x)$
- ▶ We need to measure the “closeness” of the two distributions  $\Rightarrow$  use the KL divergence

$$\text{KL}[\pi(y|x) || p(y|x, \theta)] = \mathbb{E}_{\pi(y|x)} \left[ \log \frac{\pi(y|x)}{p(y|x, \theta)} \right]$$

- ▶ We want this KL to be small for all  $x$  sampled from  $\pi(x)$

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{\pi(x)} [\text{KL}[\pi(y|x) || p(y|x, \theta)]]$$



# Probabilistic modelling: linear regression

Step 2: estimating  $\theta$  using **maximum likelihood estimation (MLE)**

Idea of MLE: for datapoints  $(x, y)$  sampled from  $\pi(x, y)$

- ▶ We want to find  $\theta^*$  such that  $p(y|x, \theta^*) \approx \pi(y|x)$
- ▶ We need to measure the “closeness” of the two distributions  $\Rightarrow$  use the KL divergence

$$\text{KL}[\pi(y|x) || p(y|x, \theta)] = \mathbb{E}_{\pi(y|x)} \left[ \log \frac{\pi(y|x)}{p(y|x, \theta)} \right]$$

- ▶ We want this KL to be small for all  $x$  sampled from  $\pi(x)$

$$\Leftrightarrow \theta^* = \arg \max_{\theta} \mathbb{E}_{\pi(x,y)} [\log p(y|x, \theta)]$$

# Probabilistic modelling: linear regression

Step 2: estimating  $\theta$  using **maximum likelihood estimation (MLE)**

Idea of MLE: for datapoints  $(\mathbf{x}, y)$  sampled from  $\pi(\mathbf{x}, y)$

- ▶ We want to find  $\theta^*$  such that  $p(y|\mathbf{x}, \theta^*) \approx \pi(y|\mathbf{x})$
- ▶ We need to measure the “closeness” of the two distributions  $\Rightarrow$  use the KL divergence

$$\text{KL}[\pi(y|\mathbf{x})||p(y|\mathbf{x}, \theta)] = \mathbb{E}_{\pi(y|\mathbf{x})} \left[ \log \frac{\pi(y|\mathbf{x})}{p(y|\mathbf{x}, \theta)} \right]$$

- ▶ We want this KL to be small for all  $\mathbf{x}$  sampled from  $\pi(\mathbf{x})$

Estimate using dataset  $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$  from  $\pi(\mathbf{x}, y)$ :

$$\theta^* = \arg \max_{\theta} \frac{1}{N} \sum_{(\mathbf{x}_n, y_n) \in \mathcal{D}} \log p(y_n | \mathbf{x}_n, \theta)$$

# Probabilistic modelling: linear regression

Step 2: estimating  $\theta$  using **maximum likelihood estimation (MLE)**

MLE: find  $\theta^*$  by

$$\theta^* = \arg \max_{\theta} \frac{1}{N} \sum_{(x_n, y_n) \in \mathcal{D}} \log p(y_n | x_n, \theta)$$

- We assumed the probabilistic model to be

$$p(y|x, \theta) = \mathcal{N}(\theta^\top \phi(x), \sigma^2)$$

# Probabilistic modelling: linear regression

Step 2: estimating  $\theta$  using **maximum likelihood estimation (MLE)**

MLE: find  $\theta^*$  by

$$\theta^* = \arg \max_{\theta} \frac{1}{N} \sum_{(x_n, y_n) \in \mathcal{D}} \log p(y_n | x_n, \theta)$$

- We assumed the probabilistic model to be

$$p(y|x, \theta) = \mathcal{N}(\theta^\top \phi(x), \sigma^2)$$

$$\Rightarrow \theta^* = \arg \max_{\theta} \frac{1}{N} \sum_{(x_n, y_n) \in \mathcal{D}} \log \mathcal{N}(\theta^\top \phi(x), \sigma^2)$$

# Probabilistic modelling: linear regression

Step 2: estimating  $\theta$  using **maximum likelihood estimation (MLE)**

MLE: find  $\theta^*$  by

$$\theta^* = \arg \max_{\theta} \frac{1}{N} \sum_{(x_n, y_n) \in \mathcal{D}} \log p(y_n | x_n, \theta)$$

- We assumed the probabilistic model to be

$$p(y|x, \theta) = \mathcal{N}(\theta^\top \phi(x), \sigma^2)$$

$$\Rightarrow \theta^* = \arg \max_{\theta} \frac{1}{N} \sum_{(x_n, y_n) \in \mathcal{D}} \log \mathcal{N}(\theta^\top \phi(x), \sigma^2)$$

$$= \arg \max_{\theta} \frac{1}{N} \sum_{(x_n, y_n) \in \mathcal{D}} -\frac{1}{2\sigma^2} \|y_n - \theta^\top \phi(x_n)\|_2^2 + \text{const}$$

# Probabilistic modelling: linear regression

Step 2: estimating  $\theta$  using **maximum likelihood estimation (MLE)**

MLE: find  $\theta^*$  by

$$\theta^* = \arg \max_{\theta} \frac{1}{N} \sum_{(x_n, y_n) \in \mathcal{D}} \log p(y_n | x_n, \theta)$$

- We assumed the probabilistic model to be

$$p(y | x, \theta) = \mathcal{N}(\theta^\top \phi(x), \sigma^2)$$

$$\begin{aligned} \Rightarrow \theta^* &= \arg \max_{\theta} \frac{1}{N} \sum_{(x_n, y_n) \in \mathcal{D}} \log \mathcal{N}(\theta^\top \phi(x), \sigma^2) \\ &= \arg \max_{\theta} \frac{1}{N} \sum_{(x_n, y_n) \in \mathcal{D}} -\frac{1}{2\sigma^2} \|y_n - \theta^\top \phi(x_n)\|_2^2 + \text{const} \\ &= \arg \min_{\theta} \frac{1}{N} \sum_{(x_n, y_n) \in \mathcal{D}} \frac{1}{2\sigma^2} \|y_n - \theta^\top \phi(x_n)\|_2^2 \end{aligned}$$

# Probabilistic modelling: linear regression

Step 2: estimating  $\theta$  using **maximum likelihood estimation (MLE)**

$$\arg \min_{\theta} \frac{1}{N} \sum_{(x_n, y_n) \in \mathcal{D}} \frac{1}{2\sigma^2} \|y_n - \theta^\top \phi(x_n)\|_2^2$$

Writing the objective in matrix form:

$$\Phi = (\phi(x_1), \dots, \phi(x_N))^\top, \mathbf{y} = (y_1, \dots, y_N)^\top$$

$$\theta^* = \arg \min_{\theta} L(\theta), \quad L(\theta) = \frac{1}{2\sigma^2} \|\mathbf{y} - \Phi\theta\|_2^2$$

► Gradient of the loss  $\nabla_{\theta} L(\theta)$ :

# Probabilistic modelling: linear regression

Step 2: estimating  $\theta$  using **maximum likelihood estimation (MLE)**

$$\arg \min_{\theta} \frac{1}{N} \sum_{(x_n, y_n) \in \mathcal{D}} \frac{1}{2\sigma^2} \|y_n - \theta^\top \phi(x_n)\|_2^2$$

Writing the objective in matrix form:

$$\Phi = (\phi(x_1), \dots, \phi(x_N))^\top, \mathbf{y} = (y_1, \dots, y_N)^\top$$

$$\theta^* = \arg \min_{\theta} L(\theta), \quad L(\theta) = \frac{1}{2\sigma^2} \|\mathbf{y} - \Phi\theta\|_2^2$$

▸ Gradient of the loss  $\nabla_{\theta} L(\theta)$ :

$$\nabla_{\theta} L(\theta) = \frac{1}{\sigma^2} \Phi^\top (\Phi\theta - \mathbf{y})$$

▸ Setting  $\nabla_{\theta} L(\theta) = 0$ :



# Probabilistic modelling: linear regression

Step 2: estimating  $\theta$  using **maximum likelihood estimation (MLE)**

$$\arg \min_{\theta} \frac{1}{N} \sum_{(x_n, y_n) \in \mathcal{D}} \frac{1}{2\sigma^2} \|y_n - \theta^\top \phi(x_n)\|_2^2$$

Writing the objective in matrix form:

$$\Phi = (\phi(x_1), \dots, \phi(x_N))^\top, \mathbf{y} = (y_1, \dots, y_N)^\top$$

$$\theta^* = \arg \min_{\theta} L(\theta), \quad L(\theta) = \frac{1}{2\sigma^2} \|\mathbf{y} - \Phi\theta\|_2^2$$

▸ Gradient of the loss  $\nabla_{\theta} L(\theta)$ :

$$\nabla_{\theta} L(\theta) = \frac{1}{\sigma^2} \Phi^\top (\Phi\theta - \mathbf{y})$$

▸ Setting  $\nabla_{\theta} L(\theta) = 0$ :

$$\Rightarrow \frac{1}{\sigma^2} \Phi^\top \Phi \theta^* = \frac{1}{\sigma^2} \Phi^\top \mathbf{y} \quad \Rightarrow \theta^* = (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}$$

# Probabilistic modelling: linear regression

Step 3: check if  $p(y|\mathbf{x}, \boldsymbol{\theta}^*)$  fits  $\pi(y|\mathbf{x})$  well

Typical approaches:

- Cross validation
- Model selection with marginal likelihood

# Probabilistic modelling: linear regression

Step 3: check if  $p(y|\mathbf{x}, \boldsymbol{\theta}^*)$  fits  $\pi(y|\mathbf{x})$  well

Typical approaches:

- Cross validation
- Model selection with marginal likelihood

If model fit is bad:

- Try another set of features  $\phi'(\mathbf{x}) \neq \phi(\mathbf{x})$
- Use other classes of models other than linear regression

# Probabilistic modelling: linear regression

Step 4: using  $p(y|x, \theta^*)$  to make predictions

Assume new test input  $x_{test}$ :

$$\theta^* = (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}$$

$$\Rightarrow p(y_{test}|x_{test}, \theta^*) = \mathcal{N}(\mathbf{y}^\top \Phi (\Phi^\top \Phi)^{-1} \phi(x_{test}), \sigma^2)$$

# Probabilistic modelling: logistic regression

Step 1: making assumptions about the output generation process

$$y = \begin{cases} 1, & \text{with probability } \rho \\ 0, & \text{with probability } 1 - \rho \end{cases}, \quad \rho = \text{sigmoid}(\boldsymbol{\theta}^\top \boldsymbol{\phi}(\boldsymbol{x}))$$

Probabilistic formulation:

- The distribution of  $y$  given  $\boldsymbol{x}$  under model assumption:

$$p(y|\boldsymbol{x}, \boldsymbol{\theta}) = \text{Bern}(\text{sigmoid}(\boldsymbol{\theta}^\top \boldsymbol{\phi}(\boldsymbol{x})))$$

# Probabilistic modelling: linear regression

Step 2: estimating  $\theta$  using **maximum likelihood estimation (MLE)**

MLE: find  $\theta^*$  by

$$\theta^* = \arg \max_{\theta} \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \log p(y|x, \theta)$$

- We assumed the probabilistic model to be

$$p(y|x, \theta) = \text{Bern}(\text{sigmoid}(\theta^\top \phi(x)))$$

# Probabilistic modelling: linear regression

Step 2: estimating  $\theta$  using **maximum likelihood estimation (MLE)**

MLE: find  $\theta^*$  by

$$\theta^* = \arg \max_{\theta} \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \log p(y|x, \theta)$$

- ▶ We assumed the probabilistic model to be

$$p(y|x, \theta) = \text{Bern}(\text{sigmoid}(\theta^\top \phi(x)))$$

$$\Rightarrow \theta^* = \arg \max_{\theta} \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \log \text{Bern}(\text{sigmoid}(\theta^\top \phi(x)))$$

# Probabilistic modelling: linear regression

Step 2: estimating  $\theta$  using **maximum likelihood estimation (MLE)**

MLE: find  $\theta^*$  by

$$\theta^* = \arg \max_{\theta} \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \log p(y|x, \theta)$$

- ▶ We assumed the probabilistic model to be

$$p(y|x, \theta) = \text{Bern}(\text{sigmoid}(\theta^\top \phi(x)))$$

$$\begin{aligned} \Rightarrow \theta^* &= \arg \max_{\theta} \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \log \text{Bern}(\text{sigmoid}(\theta^\top \phi(x))) \\ &= \arg \max_{\theta} \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} y \log \hat{y}(x; \theta) + (1 - y) \log(1 - \hat{y}(x; \theta)), \\ &\quad \hat{y}(x; \theta) = \text{sigmoid}(\theta^\top \phi(x)) \end{aligned}$$



# Probabilistic modelling: linear regression

Step 2: estimating  $\theta$  using **maximum likelihood estimation (MLE)**

$$\arg \max_{\theta} L(\theta), \quad L(\theta) = \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} y \log \hat{y}(x; \theta) + (1 - y) \log(1 - \hat{y}(x; \theta)),$$

$$\hat{y}(x; \theta) = \text{sigmoid}(\theta^\top \phi(x))$$

► Gradient of the loss  $\nabla_{\theta} L(\theta)$ :

$$\nabla_{\theta} L(\theta) = \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} [y - \hat{y}(x; \theta)] \phi(x)$$

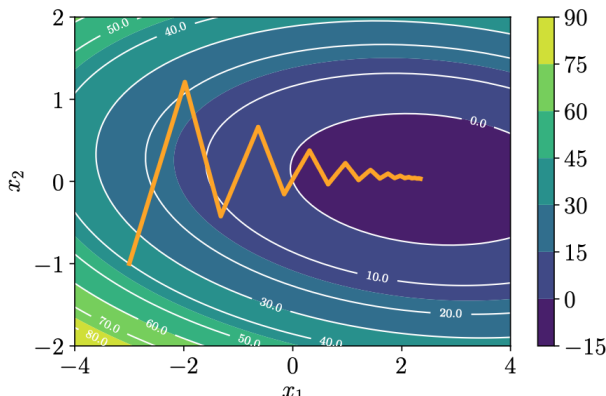
No analytic solutions!

# Gradient descent based optimisation

Algorithm: Gradient Descent (gradient **ascent** in MLE case)

Define **starting point**  $\theta_0$ , sequence of **step sizes**  $\gamma_t$ , set  $t \leftarrow 0$ .

1. Set  $\theta_{t+1} = \theta_t + \gamma_t \nabla_{\theta} L(\theta_t)$ ,  $t \leftarrow t + 1$
2. Repeat 1 until stopping criterion.



# Probabilistic modelling: logistic regression

Step 3: check if  $p(y|\mathbf{x}, \boldsymbol{\theta}^*)$  fits  $\pi(y|\mathbf{x})$  well

Typical approaches:

- Cross validation
- Model selection with marginal likelihood

If model fit is bad:

- Try another set of features  $\phi'(\mathbf{x}) \neq \phi(\mathbf{x})$
- Use other classes of models other than logistic regression

# Probabilistic modelling: logistic regression

Step 4: using  $p(y|x, \theta^*)$  to make predictions

Assume new test input  $x_{test}$ :

$\theta^*$  obtained by gradient descent

$$\Rightarrow p(y_{test}|x_{test}, \theta^*) = \text{Bern}(\text{Sigmoid}((\theta^*)^\top \phi(x_{test})))$$

# Probabilistic modelling & MLE: summary

Have you ever wondered about the following questions:

- Why using  $\ell_2$  loss in many regression problems?

**A:** We assume the model to be  $p(y|x, \theta) = \mathcal{N}(\theta^\top \phi(x), \sigma^2)$ ,  
and fit  $\theta$  using MLE

# Probabilistic modelling & MLE: summary

Have you ever wondered about the following questions:

- Why using  $\ell_2$  loss in many regression problems?  
**A:** We assume the model to be  $p(y|x, \theta) = \mathcal{N}(\theta^\top \phi(x), \sigma^2)$ ,  
and fit  $\theta$  using MLE
- Where does the cross-entropy loss come from?  
**A:** It comes from MLE, and in binary classification using  
 $p(y|x, \theta) = \text{Bern}(\text{sigmoid}(\theta^\top \phi(x)))$

# Probabilistic modelling & MLE: summary

Have you ever wondered about the following questions:

- ▶ Why using  $\ell_2$  loss in many regression problems?  
**A:** We assume the model to be  $p(y|x, \theta) = \mathcal{N}(\theta^\top \phi(x), \sigma^2)$ ,  
and fit  $\theta$  using MLE
- ▶ Where does the cross-entropy loss come from?  
**A:** It comes from MLE, and in binary classification using  
 $p(y|x, \theta) = \text{Bern}(\text{sigmoid}(\theta^\top \phi(x)))$
- ▶ What is a good principle for choosing a good loss function?  
**A:** Build a probabilistic model for the data generation process,  
and fit the parameters using MLE (or MAP, posterior inference)

# Exercises

Finish relevant exercises in the exercise sheet

- You should be able to derive MLE objectives from probabilistic model assumptions, and vice versa

Next lecture: convergence of gradient descent

Pre-requisite knowledge: Eigen-decomposition

(See e.g., <https://youtu.be/xgZ8oK9Wxzg> or search relevant videos from e.g., 3Blue1Brown)