

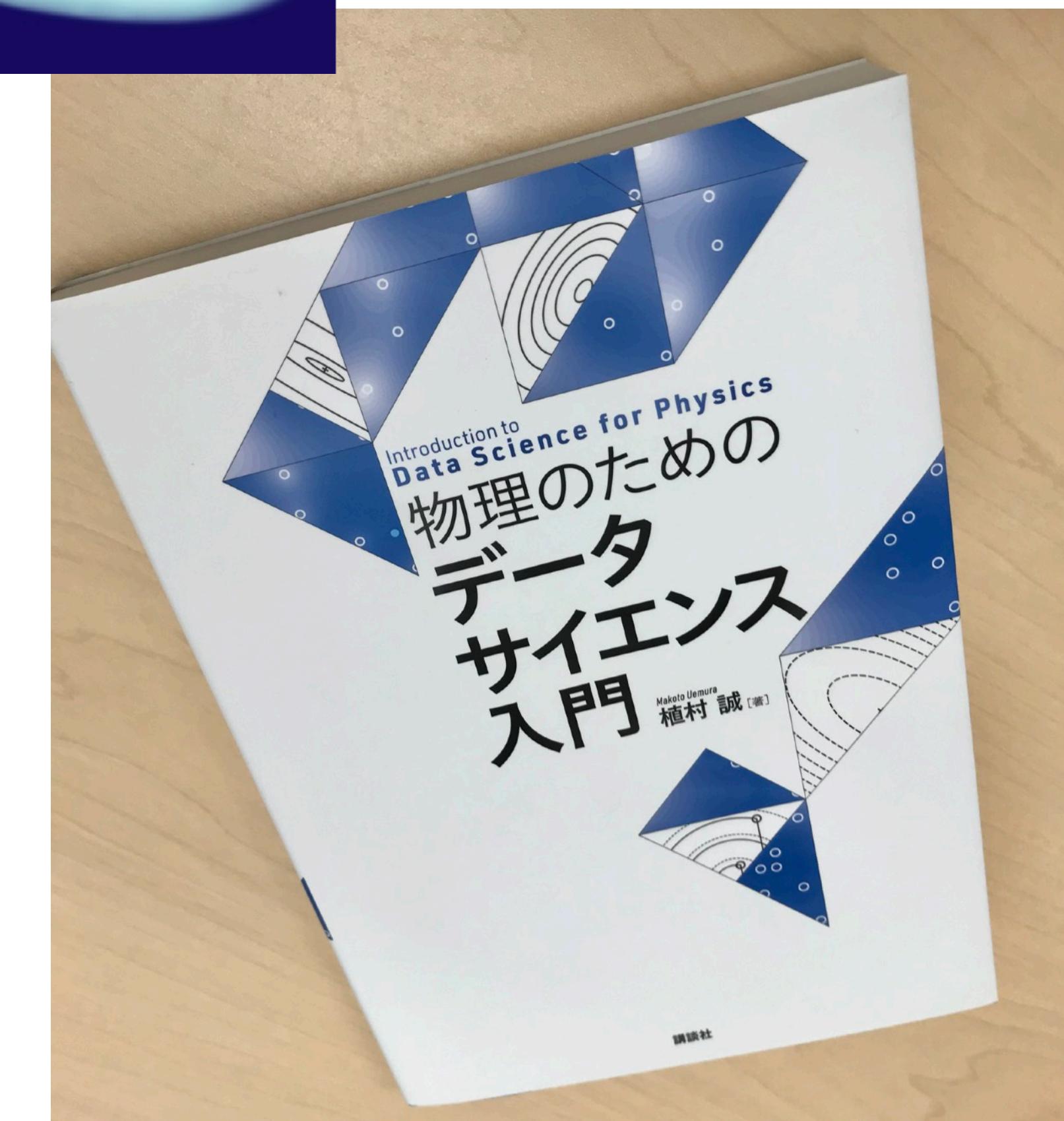
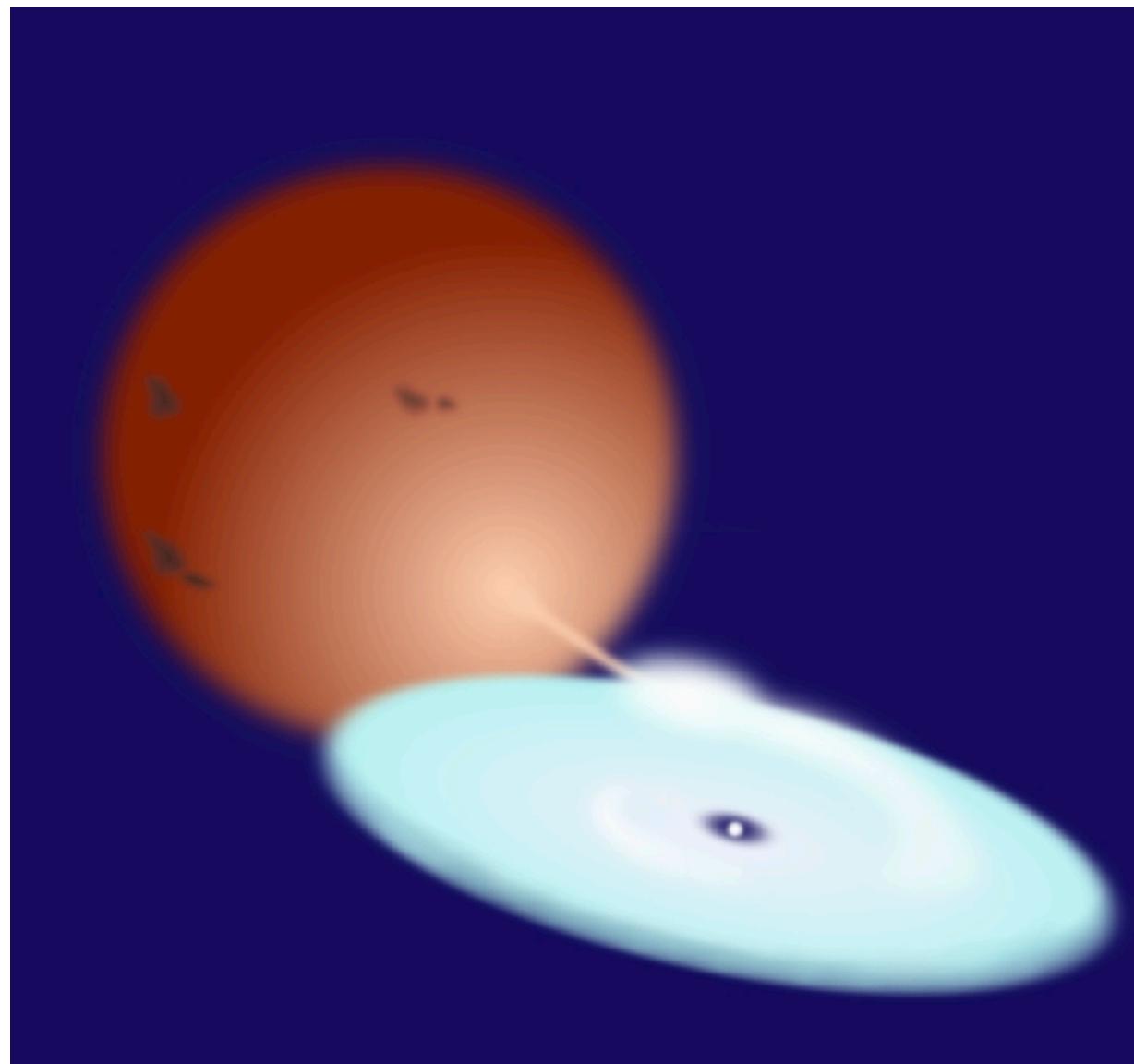
# **Statistical modeling + Shallow ML**

**AstroAI Asian Network Summer School, Sept 2-6, 2024  
DAY 1 Lecture**

**Makoto Uemura (Hiroshima University)**

# About Me

- Time Domain Astronomy (the physics of accretion & ejection around compact objects)
- Data Science for Astronomy
  - “Introduction to Data Science for Physics” 2023 Kodansha
    - Section 0: Data Science and Machine Learning, What excites you?
    - Section 1: Estimation and Test
    - Section 2: High dimensional model
    - **Section 3: Bayesian modeling**
    - Section 4: Markov-chain Monte Carlo (MCMC)
    - **Section 5: Regularization and Sparse modeling**
    - **Section 6: Classification models**
    - **Section 7: Gaussian process**
    - Section 8: Neural network
    - Appendix A: Python programs



# Contents

4 topics x (30-min lecture + 15-min hands-on exercise)

- Statistical modeling
- Regularization
- Gaussian process
- Classification models

# 1. Statistical modeling

1-1 Bayes' theorem, posterior, prior, likelihood

1-2 Generalization error, information criterion, cross-validation

# Terms and Notation

- Scalar  $x$ , Vector  $\mathbf{x}$  (bold), Matrix  $\mathbf{A} : \mathbf{y} = (y_1, y_2, \dots, y_N), \mathbf{y} = \mathbf{Ax}$
- Gaussian (normal) distribution, and its sample :
  - $p(z) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(z-\mu)^2}{2\sigma^2} \right\} = \mathcal{N}(\mu, \sigma^2), \quad z \sim \mathcal{N}(\mu, \sigma^2)$
- Conditional probability of  $y$  given  $\theta : p(y|\theta)$
- Norm of a vector :  $\|\mathbf{x}\|_p = \left( \sum_i |x_i|^p \right)^{1/p}, \|\mathbf{x}\|_2 = \sqrt{x_1 + x_2 + \dots + x_N}$
- $\mathbf{y} = f(\mathbf{x}; \theta)$ :  $\mathbf{y}$  target variables,  $\mathbf{x}$  explanatory variables (, or predictors or features),  $\theta$  model parameters

If you find any terms unclear later in this lecture, please refer back to this slide.

# 1-1 Bayes' theorem, posterior, prior, likelihood

# Bayes' Theorem and its Cast

## Posterior probability and Likelihood

$$p(\theta|y) = \frac{p(y|\theta)p(\theta)}{p(y)}$$

- Posterior probability  $p(\theta|y)$ 
  - The goal of Bayesian inference is to estimate the posterior distribution.
- Likelihood function  $p(y|\theta)$ 
  - joint probability of observing the data  $y$  given the parameters  $\theta$
  - In frequentist statistics, the best model is that maximizes the likelihood.

# Bayes' Theorem and its Cast

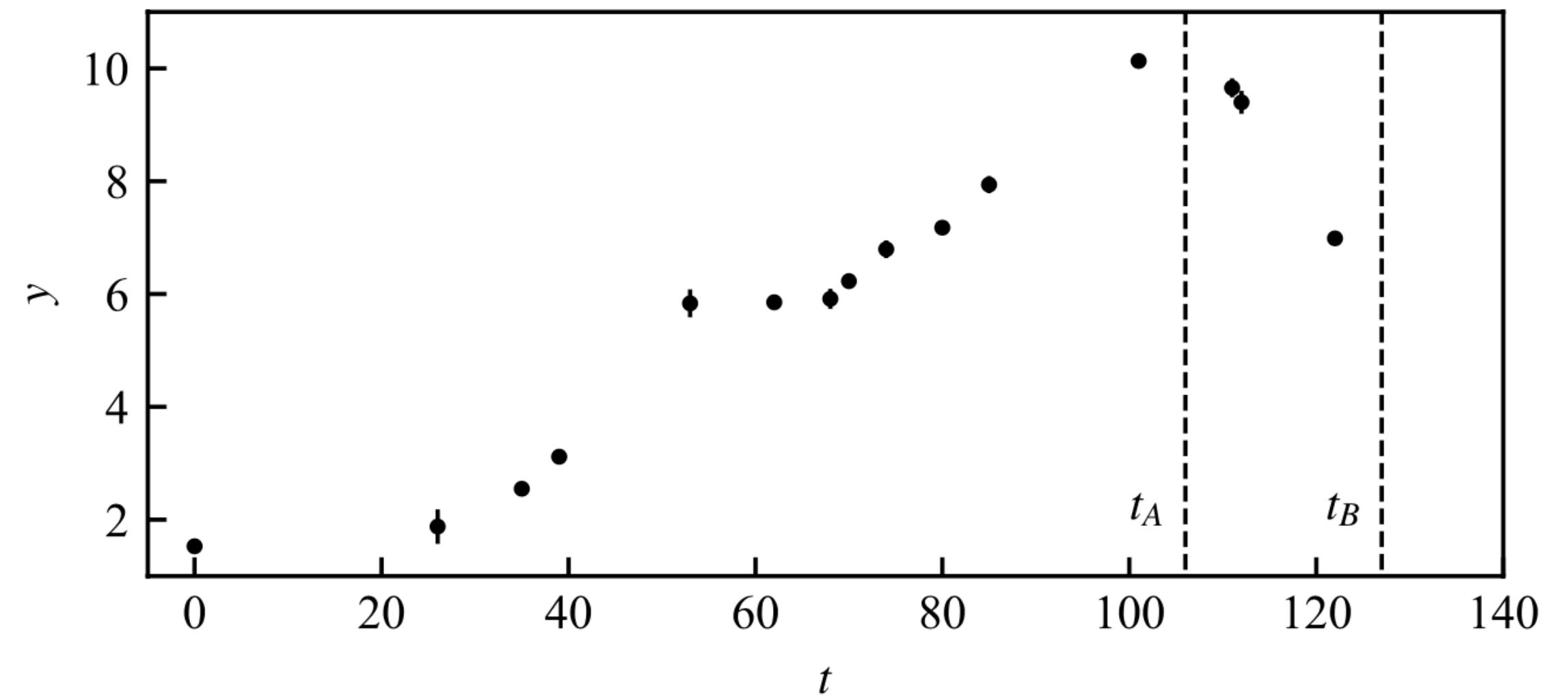
## Prior probability

$$p(\theta|y) = \frac{p(y|\theta)p(\theta)}{p(y)}$$

- Prior probability  $p(\theta)$ 
  - A key component of Bayesian models
  - allows us to incorporate existing knowledge or assumptions about the parameters into the model.
  - helps to articulate relationships between parameters more clearly.

# Ex. Smooth curve

- Time-series data (15 data points in  $t=0-140$ )
- Equally-spaced 141  $\mu$
- Current slope is similar to the last slope
- Inputs: 15 Data  $(t_i, y_i)$
- Parameters: 141  $\mu_i$  &  $\lambda$
- Estimating the posterior of the parameters with MCMC ( $\lambda$  is given by hand)

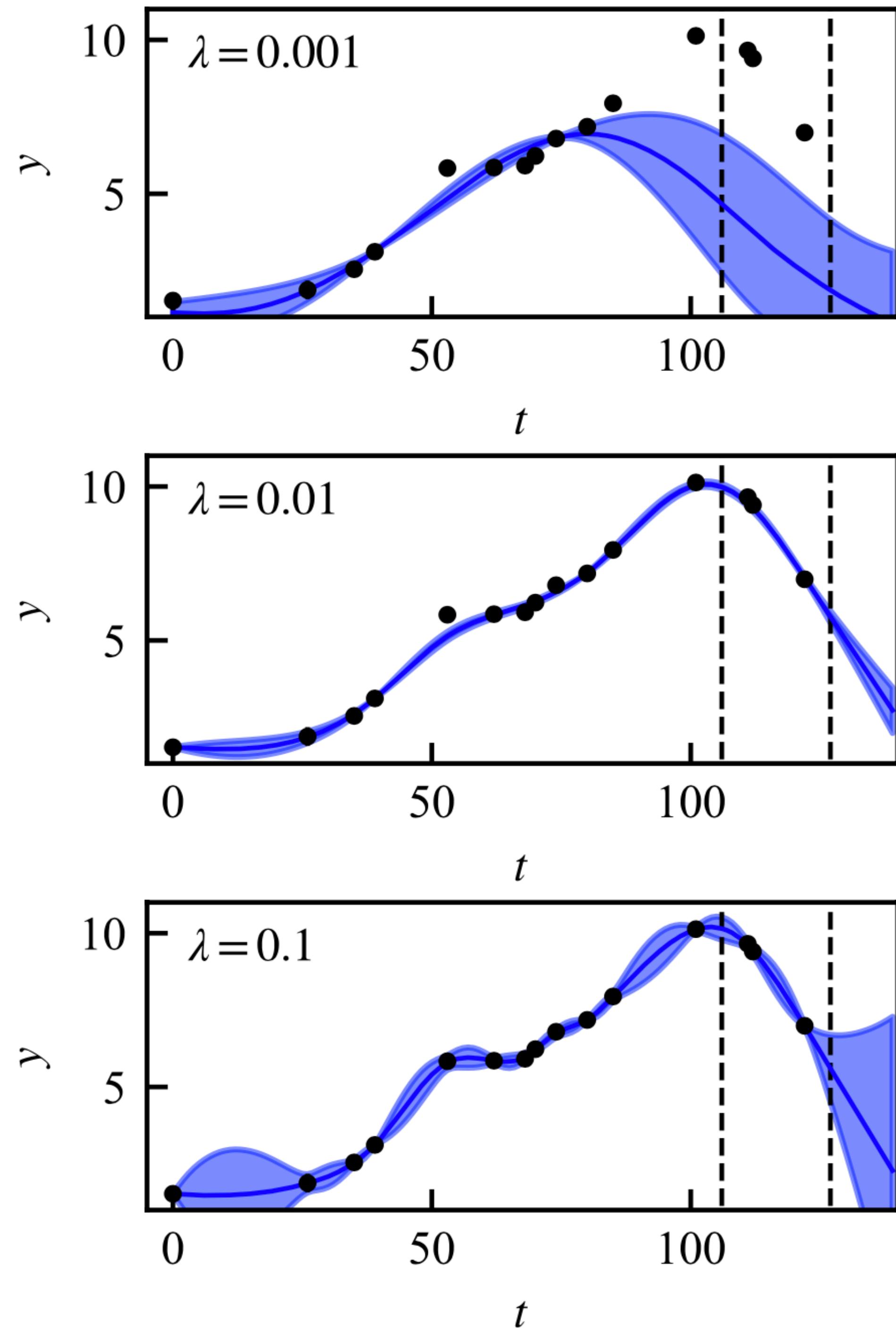


$$\begin{aligned}
 p(y_i|\mu_j) &= \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp \left\{ -\frac{(y_i - \mu_j)^2}{2\sigma_i^2} \right\} \quad (\text{ただし、 } t_i = t_j) \\
 p(\mu_j) &= \frac{1}{\sqrt{2\pi\lambda^2}} \exp \left\{ -\frac{((\mu_j - \mu_{j-1}) - (\mu_{j-1} - \mu_{j-2}))^2}{2\lambda^2} \right\} \\
 &= \frac{1}{\sqrt{2\pi\lambda^2}} \exp \left\{ -\frac{(\mu_j - 2\mu_{j-1} + \mu_{j-2})^2}{2\lambda^2} \right\}
 \end{aligned}$$

# Ex. Smooth curve: Result

$$\begin{aligned}
 p(\mu_j) &= \frac{1}{\sqrt{2\pi\lambda^2}} \exp \left\{ -\frac{((\mu_j - \mu_{j-1}) - (\mu_{j-1} - \mu_{j-2}))^2}{2\lambda^2} \right\} \\
 &= \frac{1}{\sqrt{2\pi\lambda^2}} \exp \left\{ -\frac{(\mu_j - 2\mu_{j-1} + \mu_{j-2})^2}{2\lambda^2} \right\}
 \end{aligned}$$

- For small  $\lambda$ 
  - Too rigid & failing to capture the underlying trend in the data.
- For large  $\lambda$ 
  - Too flexible  $\rightarrow$  Overfit (see the next section)
- For moderate  $\lambda$ 
  - Smooth and capturing the trend



# Ex. Smooth curve: lessons learned

$$p(\mu_j) = \frac{1}{\sqrt{2\pi\lambda^2}} \exp \left\{ -\frac{(\mu_j - \mu_{j-1})^2}{2\lambda^2} \right\}$$

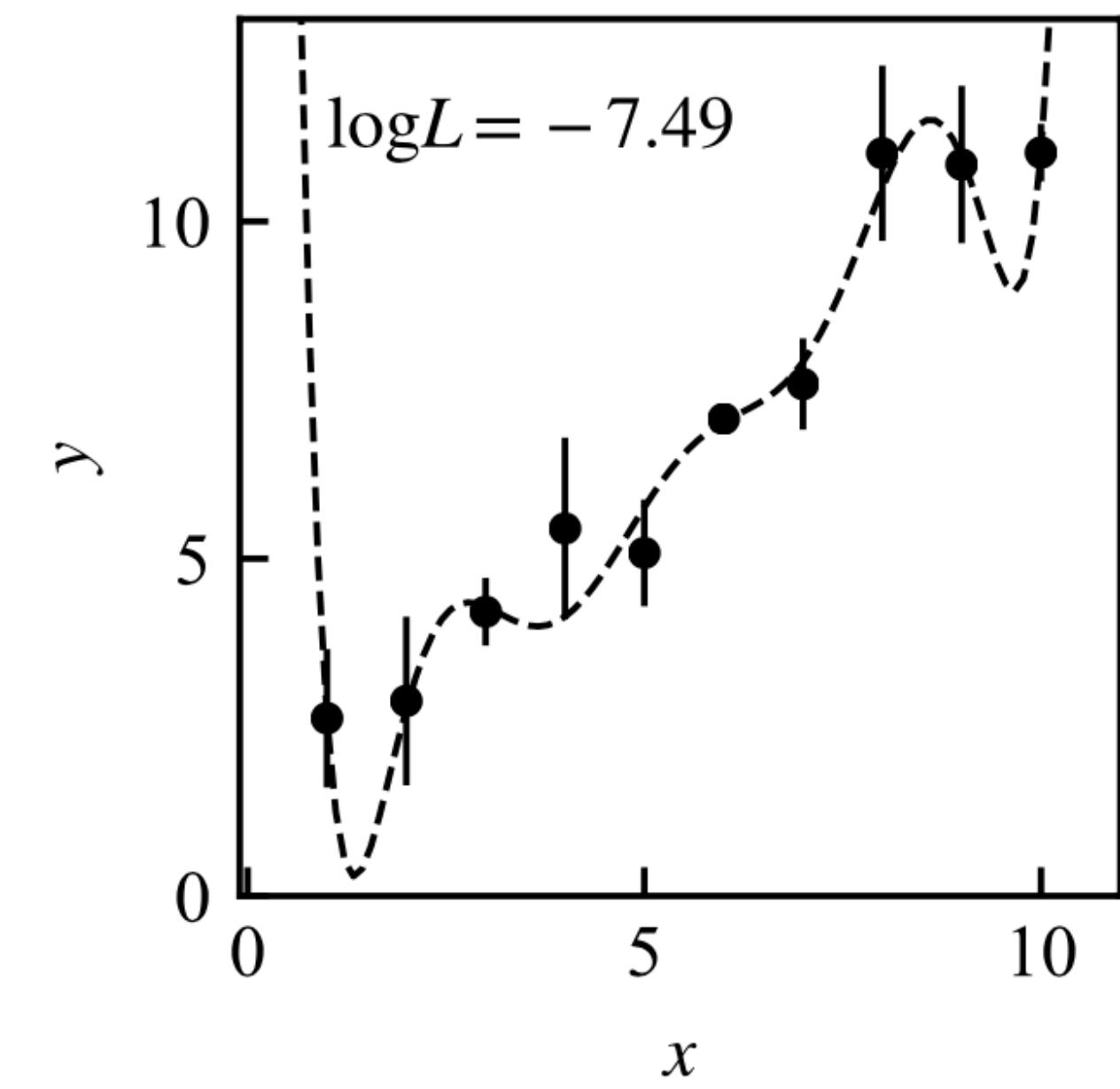
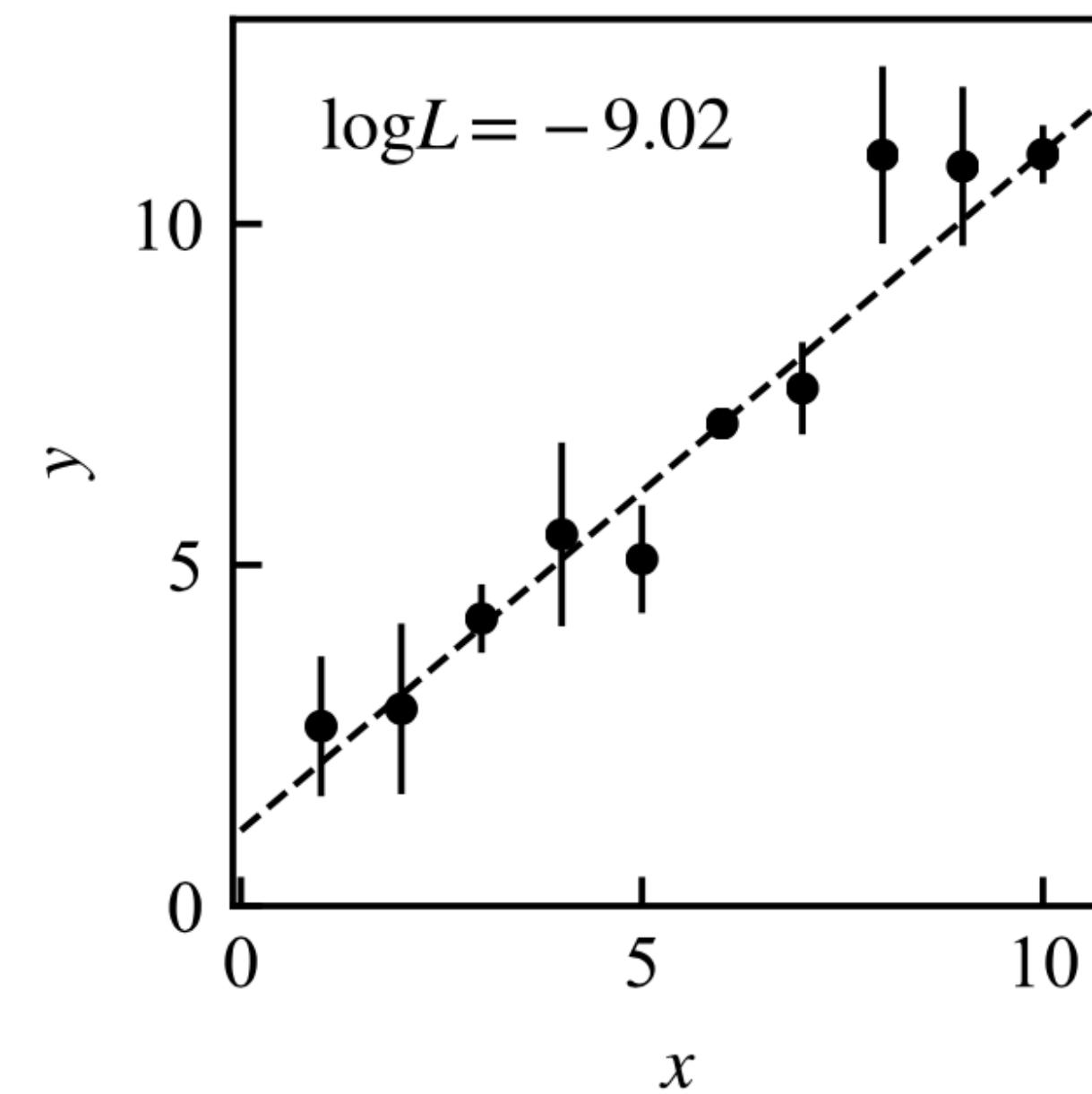
$$\begin{aligned} p(\mu_j) &= \frac{1}{\sqrt{2\pi\lambda^2}} \exp \left\{ -\frac{((\mu_j - \mu_{j-1}) - (\mu_{j-1} - \mu_{j-2}))^2}{2\lambda^2} \right\} \\ &= \frac{1}{\sqrt{2\pi\lambda^2}} \exp \left\{ -\frac{(\mu_j - 2\mu_{j-1} + \mu_{j-2})^2}{2\lambda^2} \right\} \end{aligned}$$

- Expressing complex relationships between the model parameters through the prior.
- not lines or polynomials, but more nuanced trends.
- Local trend model
  - It can be extended to
    - Auto-regressive model = the coefficients of  $\mu_{j-1}, \mu_{j-2}$  are also estimated from the data.
    - State-space model

# **1-2 Generalization error, information criteria, cross-validation**

# Polynomial Regression, as an example

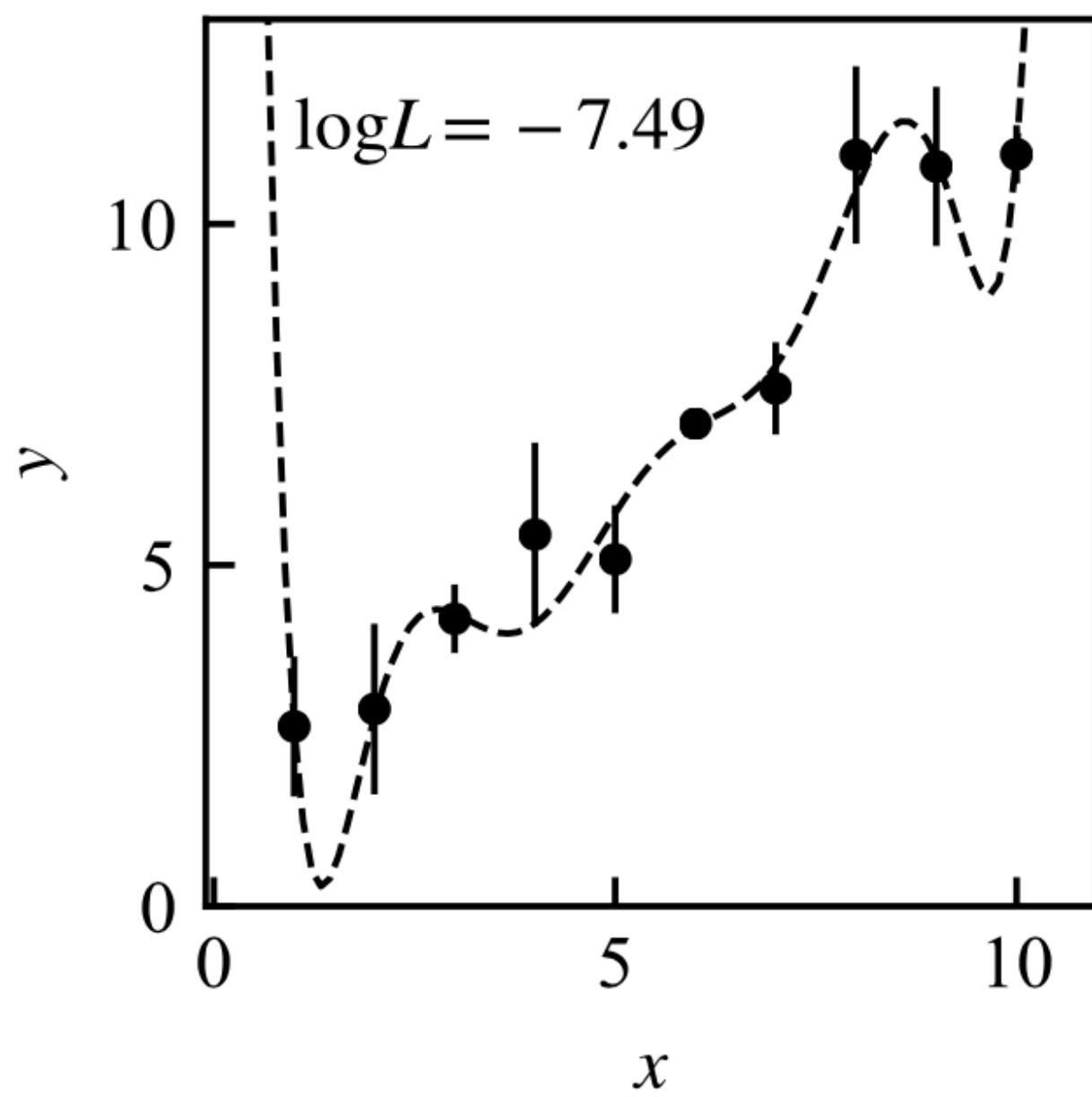
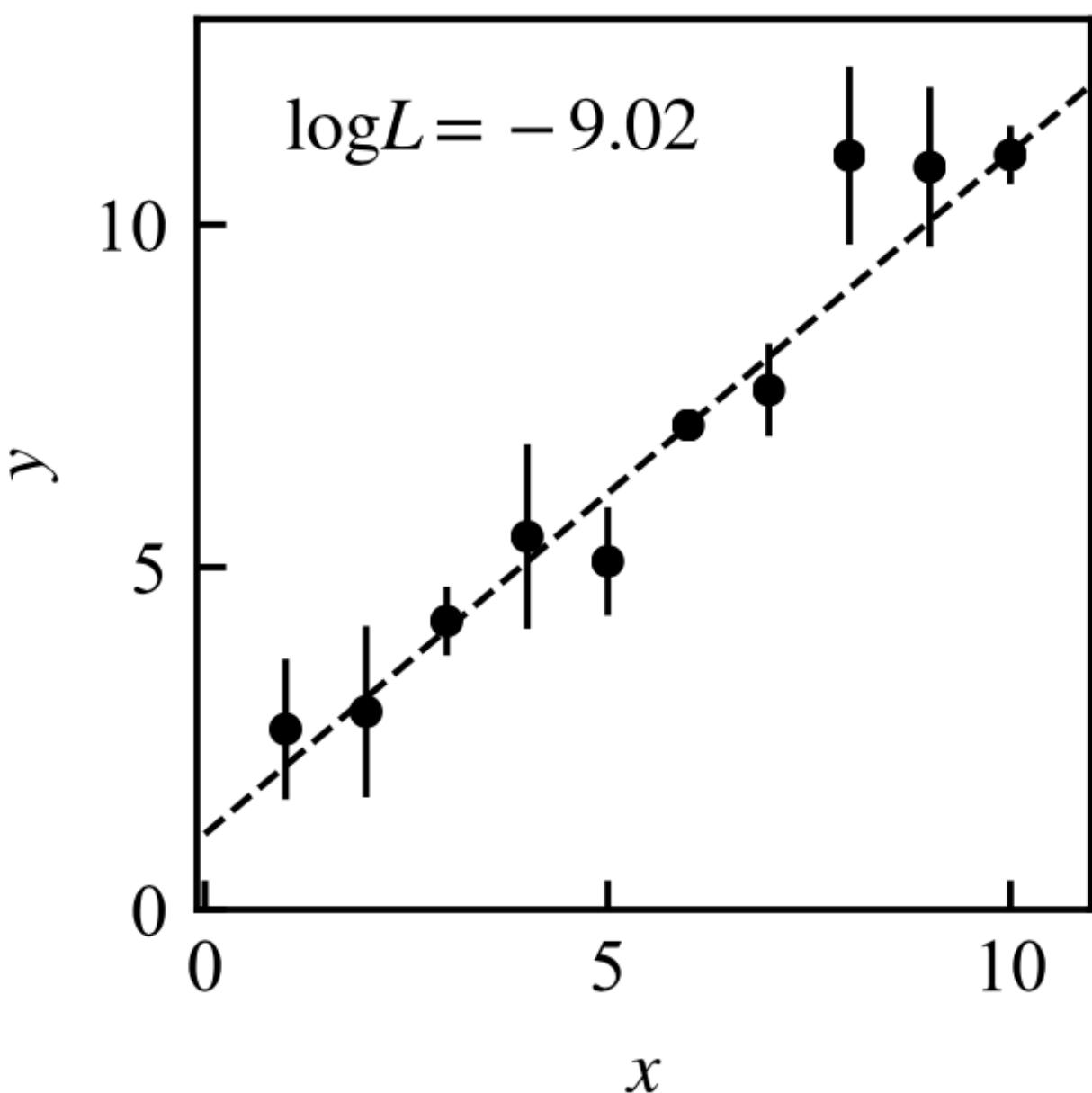
- Left panel: 1st degree = line
- Right panel: 8th degree polynomial
- The eighth-degree polynomial has a higher (log-)likelihood.
- Do you prefer the 8th order model?



# Overfit and model selection

- Overfit = a statistical model describes random noise
  - The likelihood of a model often increases with the number of model parameters.
  - In high-dimensional problems, evaluating a model solely based on likelihood becomes inappropriate.
- Model selection:
  - If the model is too simple → it cannot adequately explain the data.
  - If the model is too complex → overfitting occurs.
  - It's essential to choose a model with the right level of complexity. = model selection
- Two approaches: Information criteria & cross-validation.

Overfit =  
Low generalization error  
↓

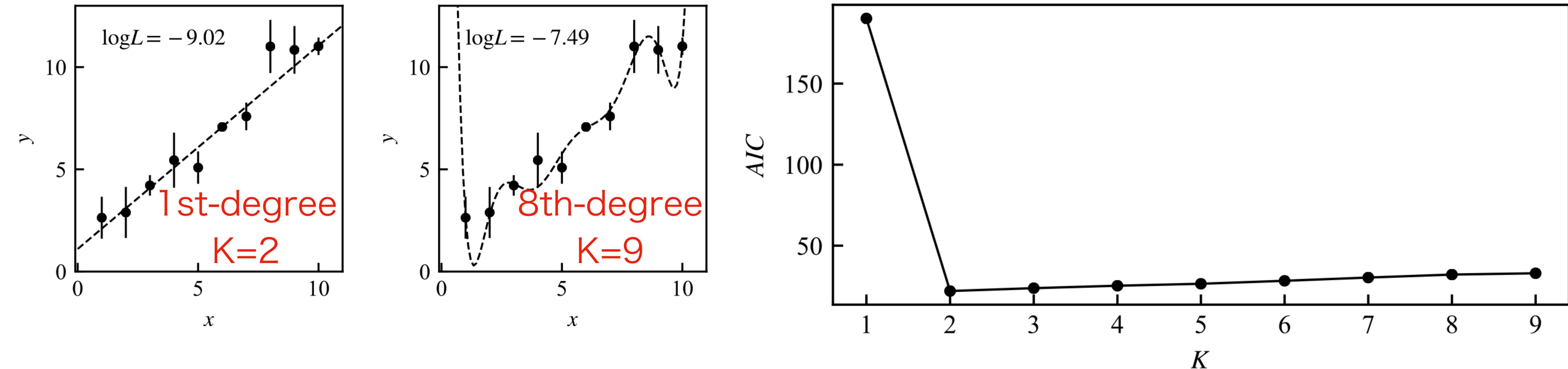


# AIC (Akaike information criterion)

$$\text{AIC} = -2(\log L(\hat{\theta}) - K) = -2 \log L(\hat{\theta}) + 2K$$

- The log-likelihood  $\log L(\hat{\theta})$  of the data can become biased when the model is too complex.
- To correct this bias, AIC introduces a penalty term, the number of model parameters  $K$ .
- Even if two models have the same likelihood, the model with fewer parameters — resulting in a smaller  $K$  — will have a lower AIC, and be considered the better model.
- Focusing not just on how well a model fits the current data, but also on how well it can predict new, unseen data. = model **predictive performance** or **generalization performance**.

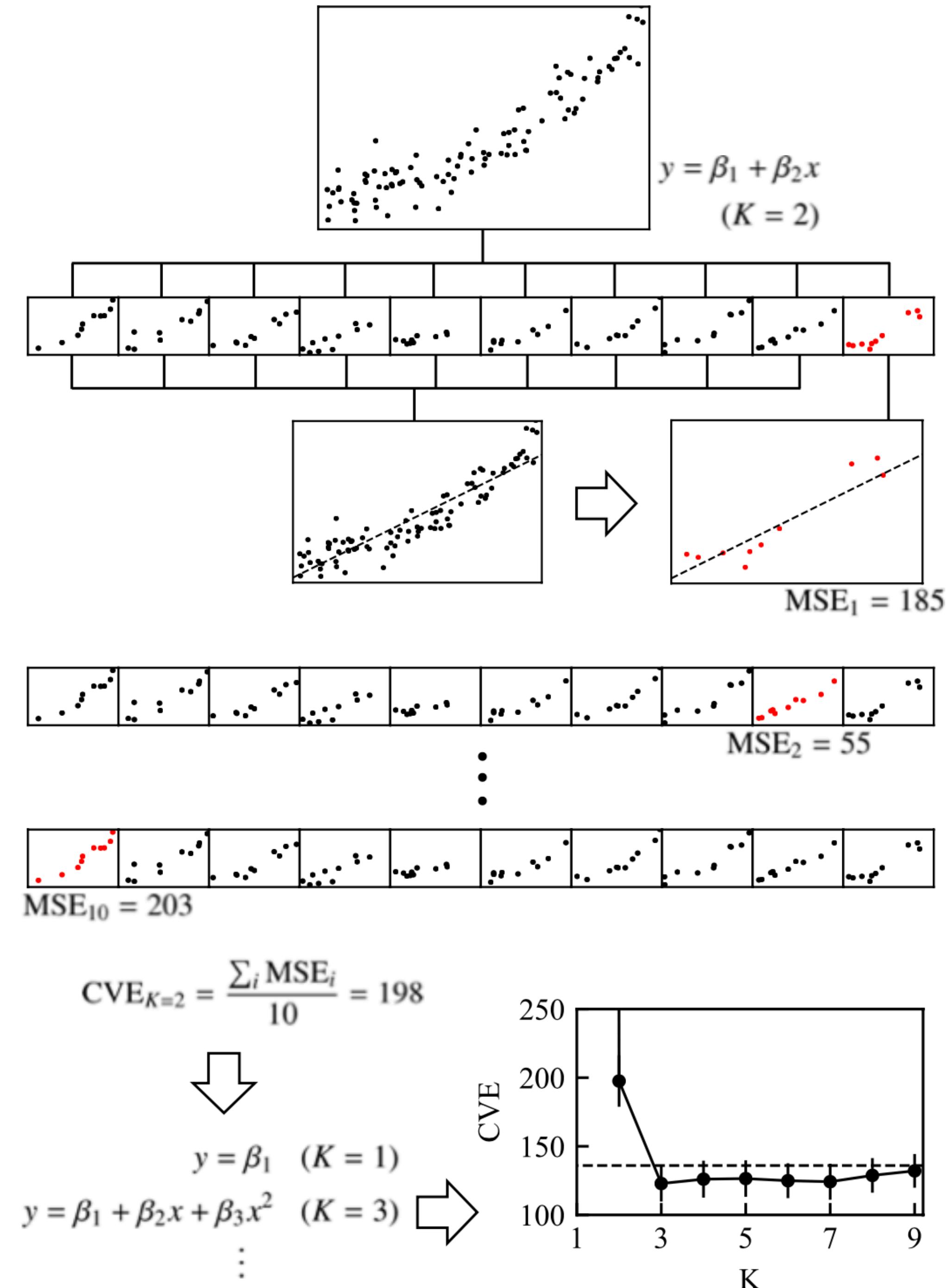
# AIC model selection for the polynomial fitting



- Models are created → the model parameters are determined using maximum likelihood estimation → the likelihood is obtained → AIC can be calculated based on the number of model parameters ( $K$ )
- When  $K = 2$ , which corresponds to a line, the AIC is minimized.
- The  $K=2$  model provides the best balance between fitting the data and avoiding unnecessary complexity.

# Cross-Validation (CV)

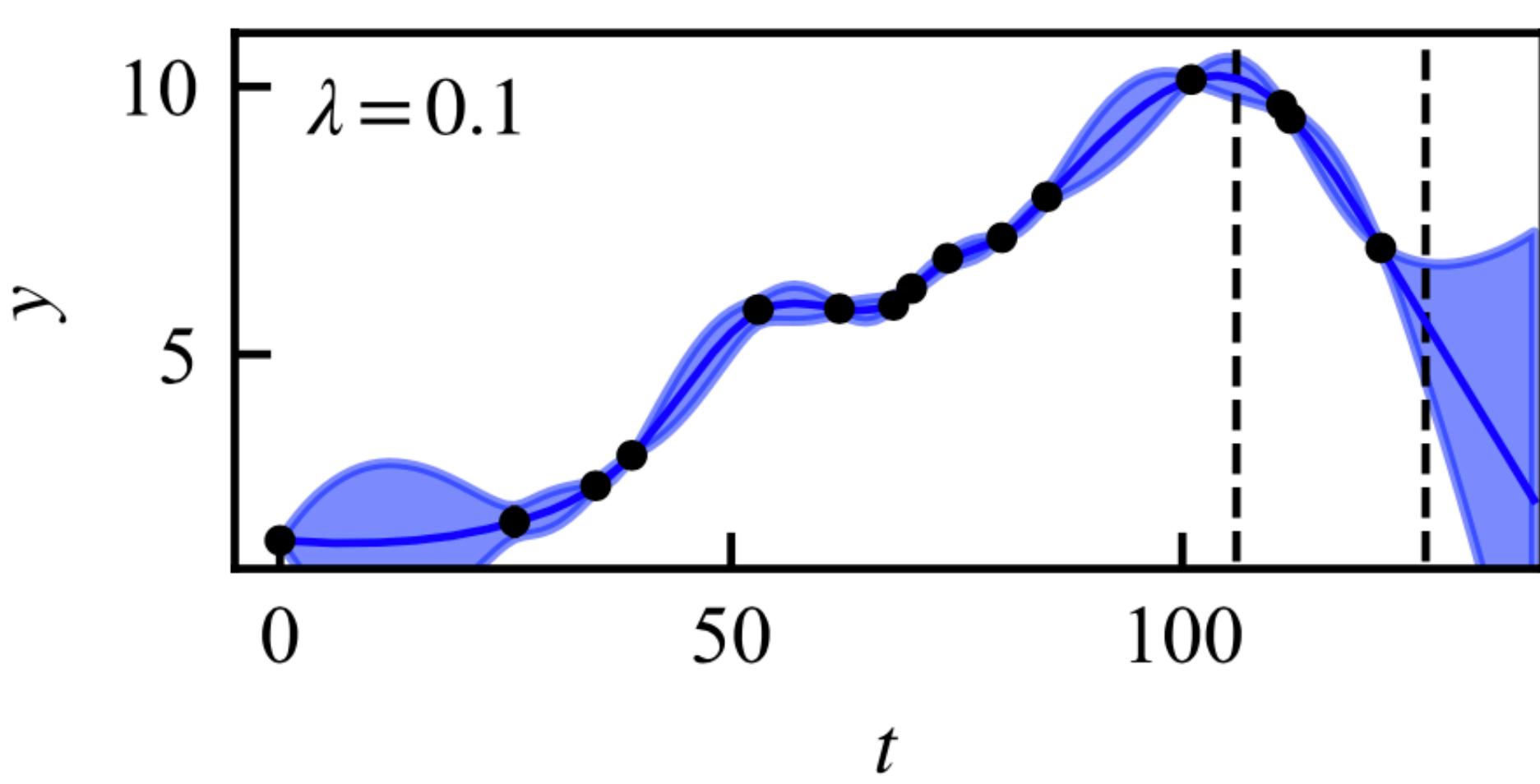
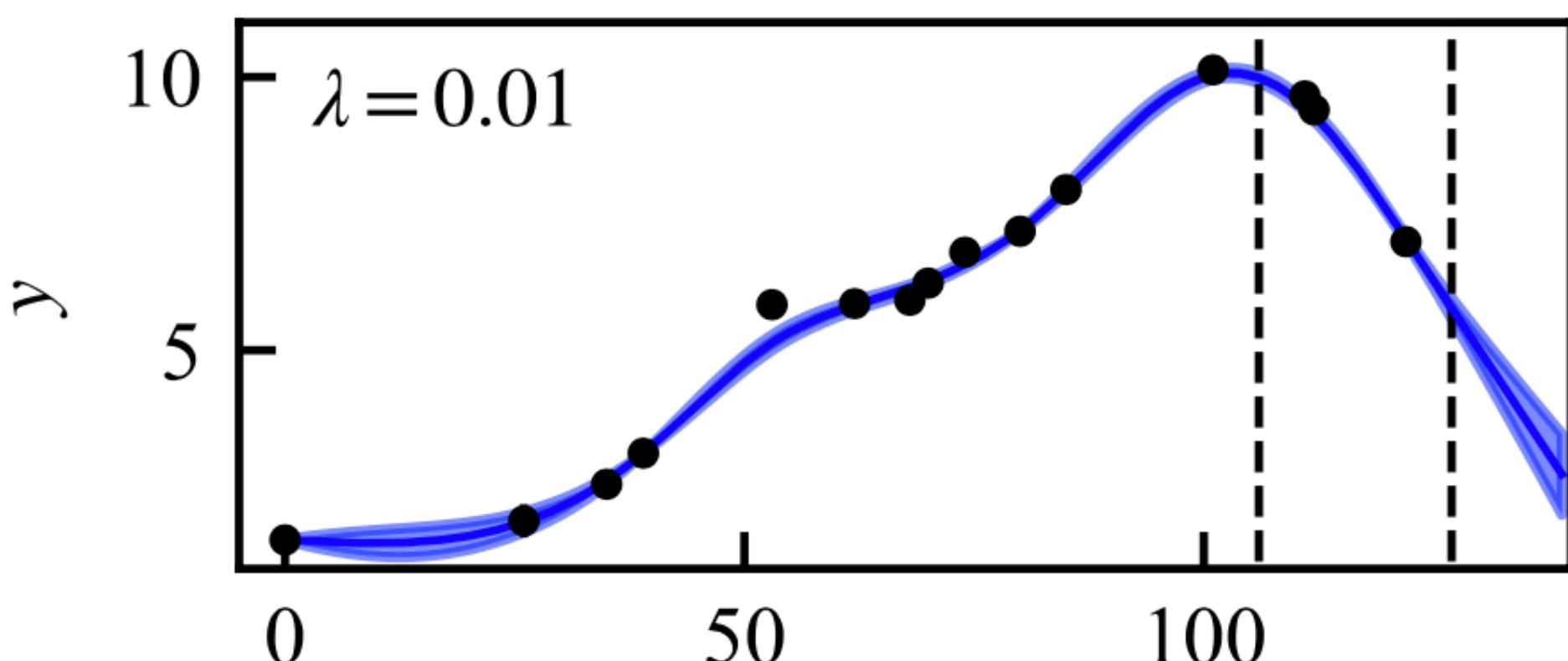
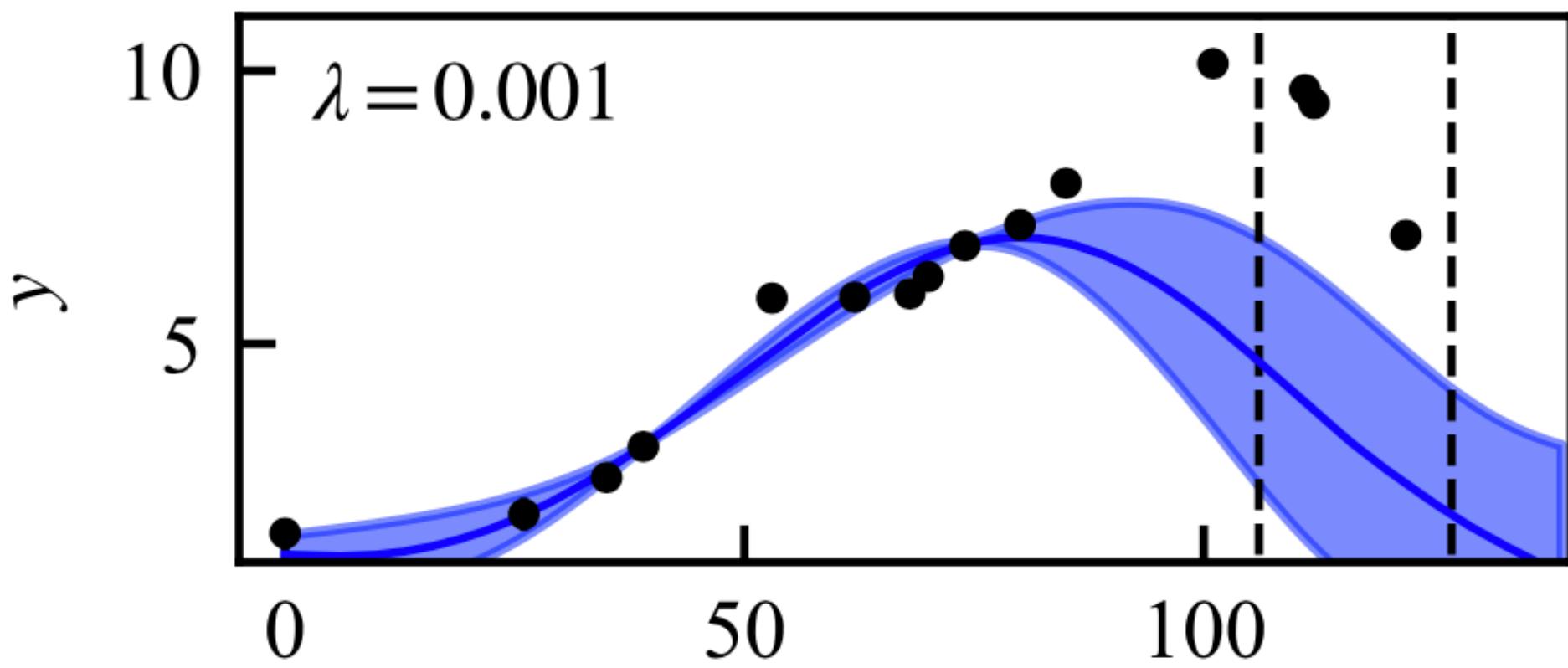
- K-fold cross-validation (right figure)
  - Divide the data into K equally-sized subsets (10 folds is common)
  - Use one subset  $k$  for validation and the remaining subsets for training
  - Obtain K metrics (such as squared errors or mean squared errors (MSE))
  - Calculate the average and standard error of the cross-validation errors (CVE)



# Hands-on exercise #1

Local trend model

[Lecture\\_Day1\\_Uemura/01\\_Bayes.ipynb](#)

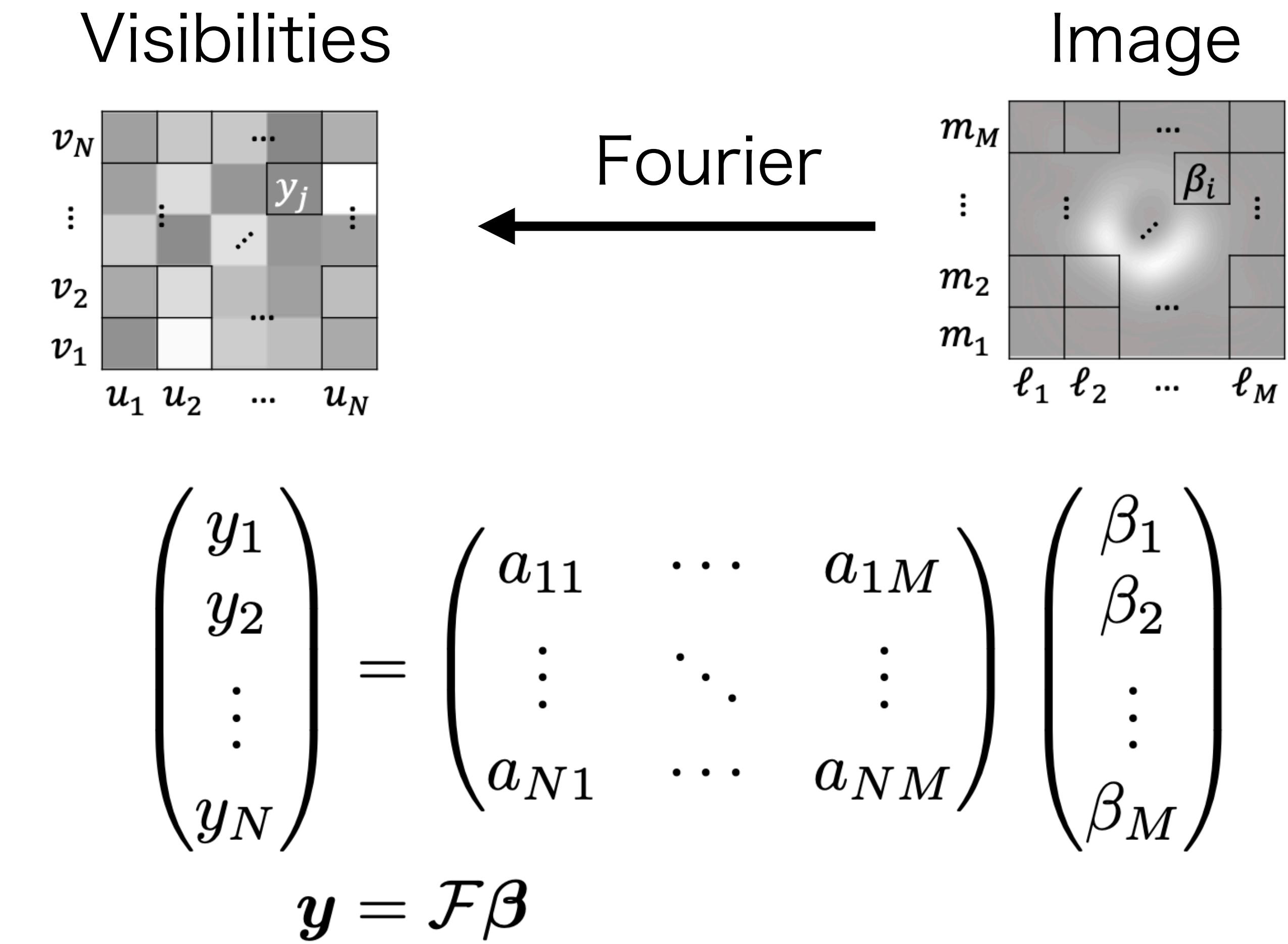


## 2. Regularization

**Key words:** LASSO ( $\ell_1$ -penalty), sparsity

# Image reconstruction in VLBI

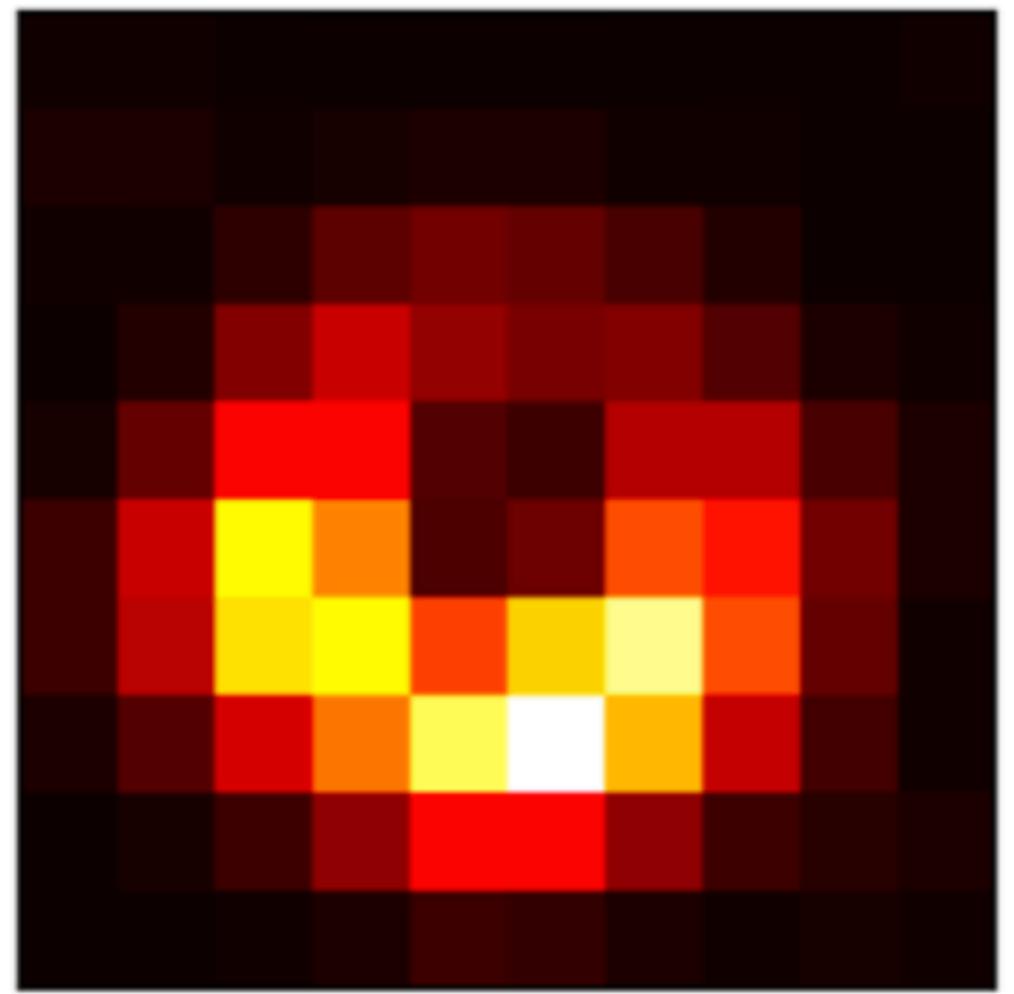
- Data: Visibilities (complex numbers)
- Estimation: Brightness distribution on the celestial sphere = 2D image
- Matrix = Fourier transform
- Number of data points  $N <$  Number of coefficients  $M$ 
  - Underdetermined system
  - = ill-posed problem



$$a_{ij} = \exp\{-2\pi i(u_i\ell_j + v_im_j)\}$$

# In maximum likelihood estimation, the solution is indeterminate

True image

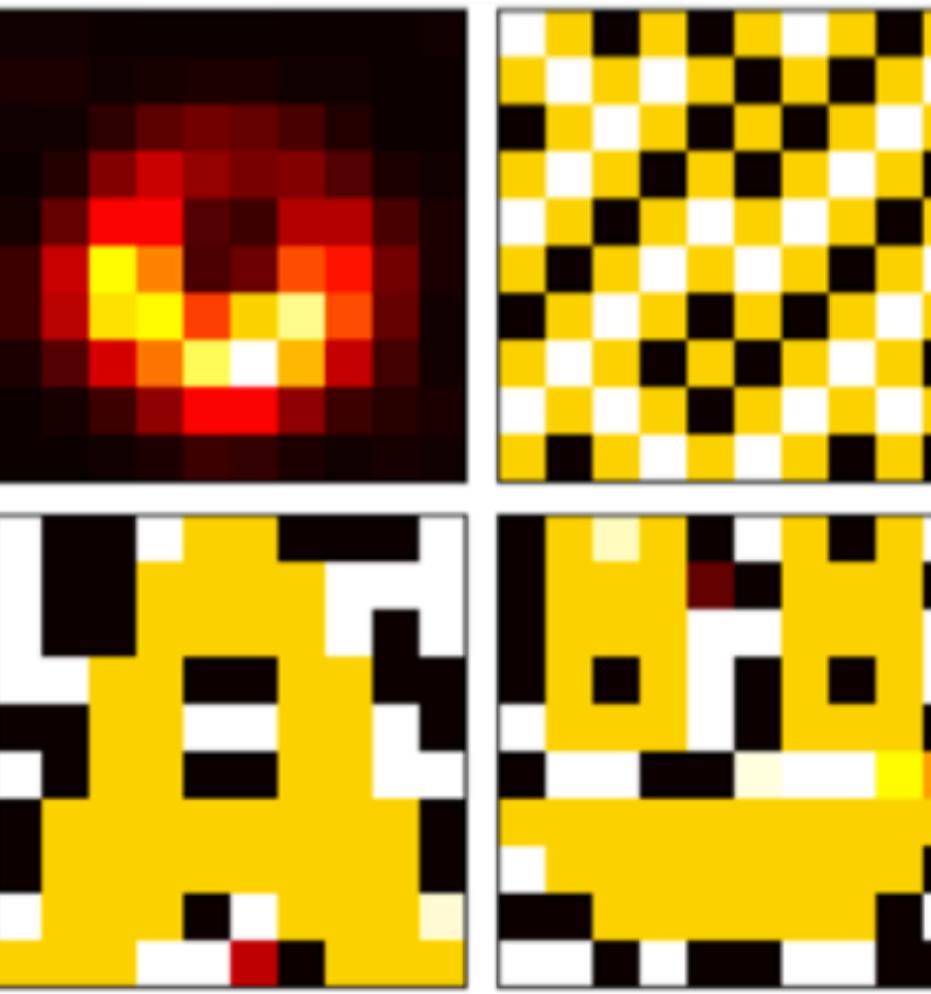


観測

Data  
(N=25)

推定

Least squares solutions



- Estimate the brightness of 100 pixels (left figure)
- Data = 25 complex visibilities = 50 pieces of information (real part + imaginary part or amplitude + phase)
- Results (right figure)
- All four images are “least squares solutions” = perfectly match the data

- So, what should we do next?
  - Does it really make sense to claim that any of the other three are “astronomical images”? → No !
  - Why do many people consider the other three to be impossible as astronomical images? → Because they have prior knowledge about what astronomical images should look like.
- Can we use that prior knowledge as a Bayesian prior probability to fill in the gaps left by the missing information?

# Least Squares Method and Regularization

## Basic terms

• Least Squares Method :  $\min_{\theta} E(\theta)$

- The problem of finding the parameter  $\theta$  that minimizes the error function,  $E(\theta) = \|y - f(\theta)\|_2^2$

• Regularized (Penalty) Least Squares Method :

$$\min_{\theta} E(\theta) + \lambda \Phi(\theta)$$

- $\lambda \Phi(\theta)$  is the regularization term

- We aim to minimize this combined function.

- $\lambda$  is the regularization coefficient.

- When  $\lambda$  is large, the method prioritizes minimizing the regularization term over fitting the data.

- When  $\lambda$  is small, the regularization term has less influence, allowing the model to fit the data more closely.

# Regularization and Bayesian Methods

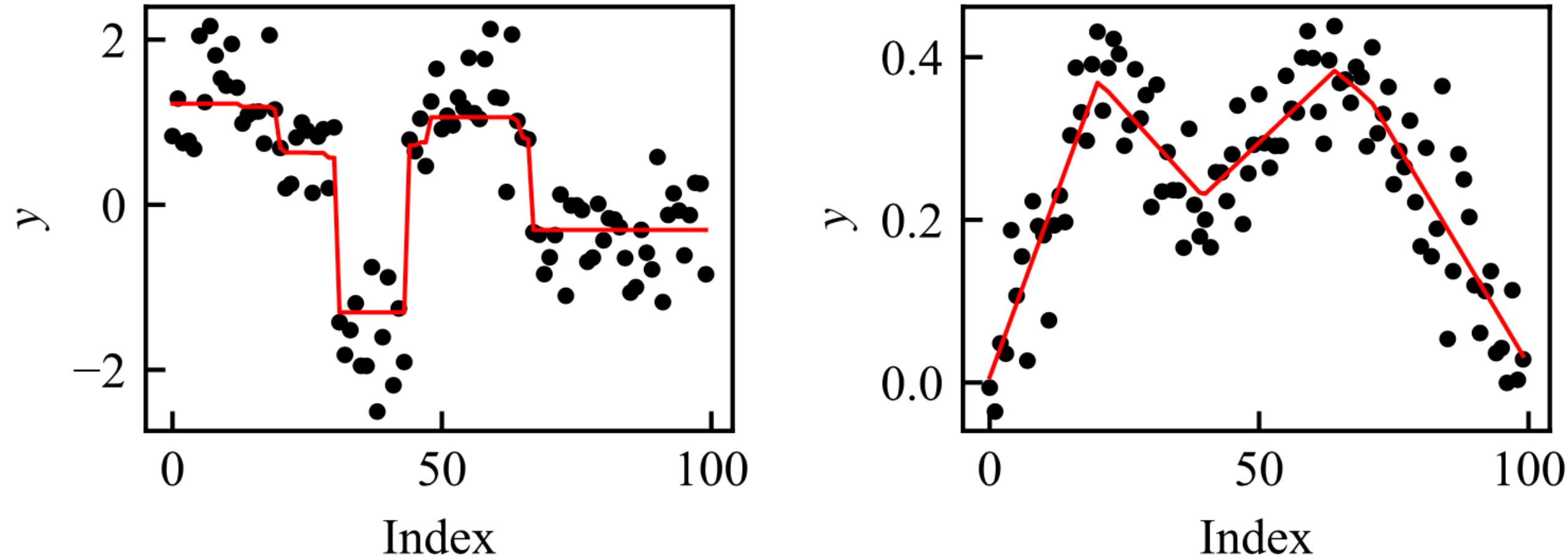
- If we think of the objective function in regularized least squares as the negative log-posterior probability:

$$-\log p(\beta | \mathbf{y}) = \frac{1}{2} \sum_i \frac{(y_i - \mathbf{x}_i^T \beta)^2}{\sigma_i^2} + \lambda \Phi(\beta) \quad \rightarrow \quad p(\beta | \mathbf{y}) \propto \prod_i \exp \left\{ -\frac{(y_i - \mathbf{x}_i^T \beta)^2}{2\sigma_i^2} \right\} \exp \{-\lambda \Phi(\beta)\}$$
$$\propto p(\mathbf{y} | \beta) p(\beta)$$

- The regularization term corresponds to the prior distribution.
- The form of the regularization term reflects our prior knowledge about the data.

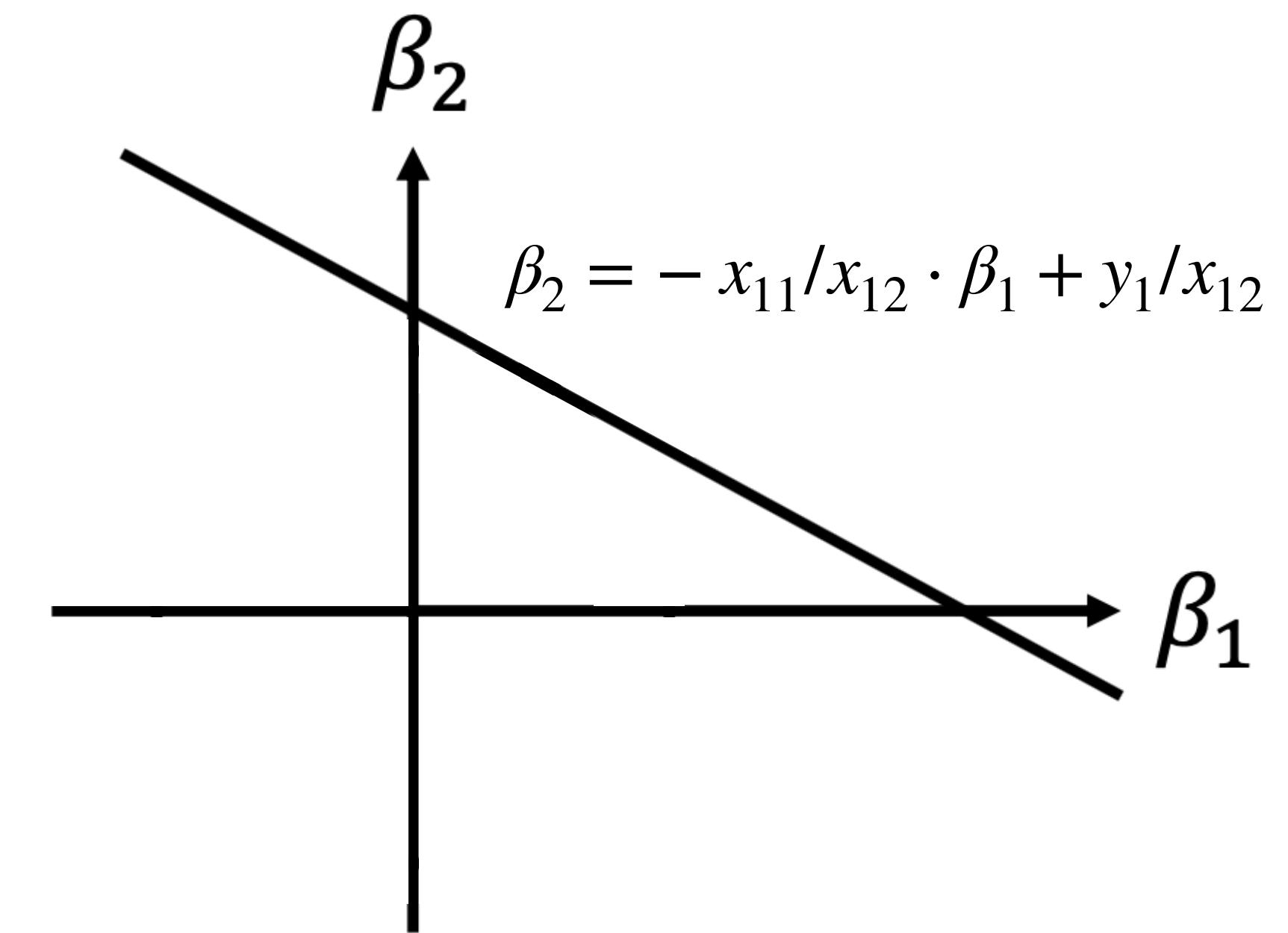
# Various regularizers

- Ridge regression :  $\lambda \|\beta\|_2^2$
- LASSO regression :  $\lambda \|\beta\|_1 \rightarrow$  sparse solution (next slides)
- Elastic net :  $\lambda(\alpha \|\beta\|_1 + (1 - \alpha) \|\beta\|_2^2)$ 
  - By adjusting the two regularization parameters  $\lambda$  and  $\alpha$ , Elastic net can sometimes yield better results than LASSO.
- Total variation minimization :  $\lambda \sum_i |\beta_{i+1} - \beta_i|$ 
  - $\ell_1$  norm in differential space.
- $\ell_1$  trend filter :  $\hat{\beta} = \arg \min_{\beta} \|\mathbf{y} - \beta\|_2^2 + \lambda \sum_i |\beta_{i+2} - 2\beta_{i+1} + \beta_i|$ 
  - sparse in second-order difference
  - Represent the data with the same line as much as possible, estimating points of trend changes.
- Group LASSO :  $\lambda \sum_k \|\beta^{(k)}\|_2$ 
$$\{\beta^{(1)}, \beta^{(2)}, \dots, \beta^{(K)}\}$$
- Maximum Entropy Method (MEM) :  $-\lambda \sum_i \beta_i \log \beta_i$ 
  - Classical regularizer often used in physics.



# LASSO regression: A simple example of undetermined problems

- Consider a low-dimensional problem that is easy to visualize
- Two model parameters  $(\beta_1, \beta_2)$
- One data point = One equation  $y_1 = \beta_1 x_{11} + \beta_2 x_{12}$
- Estimate  $(\beta_1, \beta_2)$
- All points on the line,  $\beta_2 = -x_{11}/x_{12} \cdot \beta_1 + y_1/x_{12}$ , represent valid solutions = Underdetermined



# LASSO regression (Least Absolute Shrinkage and Selection Operator; Tibshirani 1996)

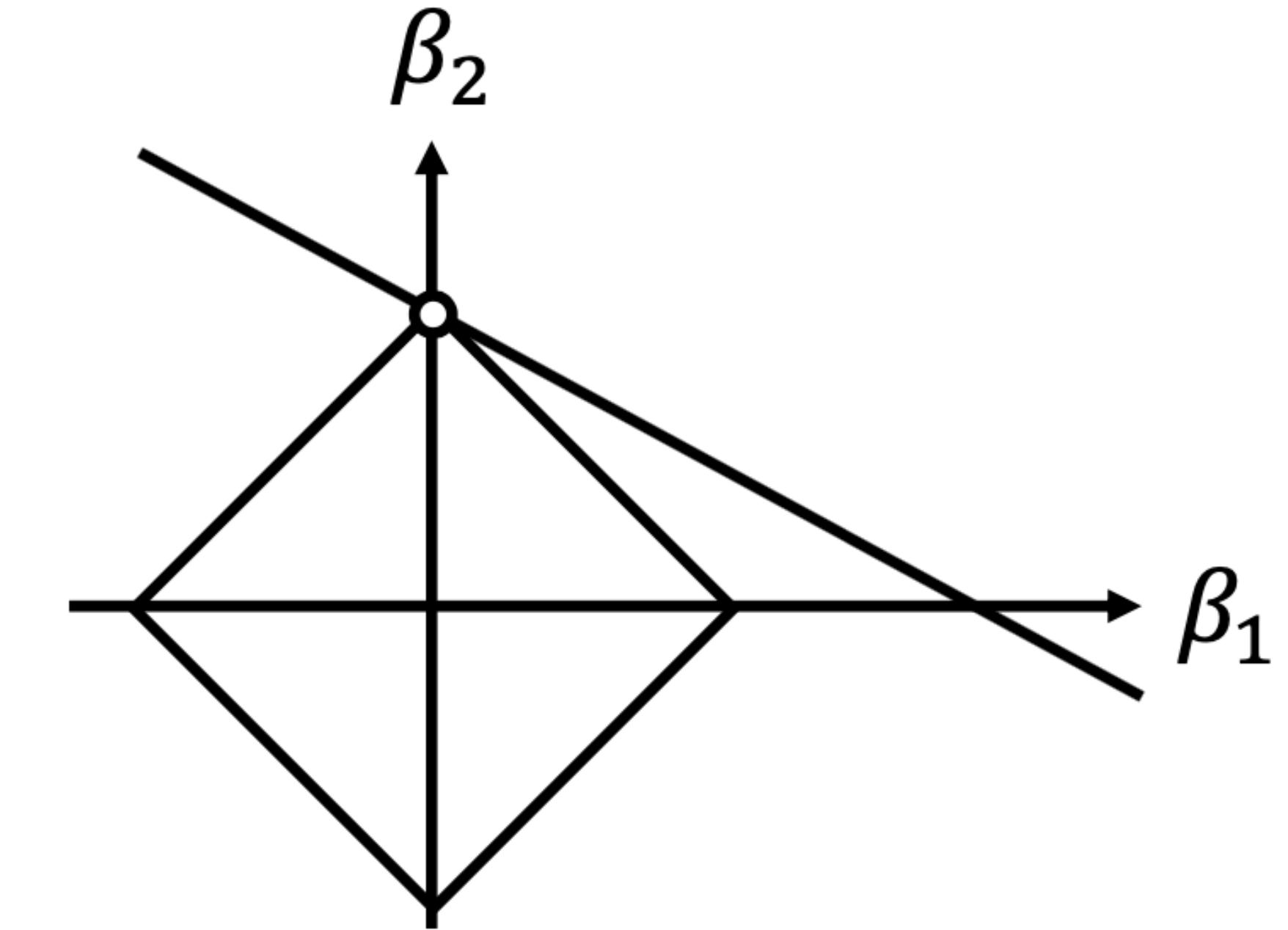
- Enforces sparsity in the coefficient vector  $\beta$  by using the L1 norm for regularization

$$\hat{\beta} = \arg \min_{\beta} \{ \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \}$$

- Constant L1 norm:

$\|\beta\|_1 = |\beta_1| + |\beta_2| + \dots + |\beta_M| = \text{const.}$  is a “diamond” shape in 2D and an octahedron in 3D

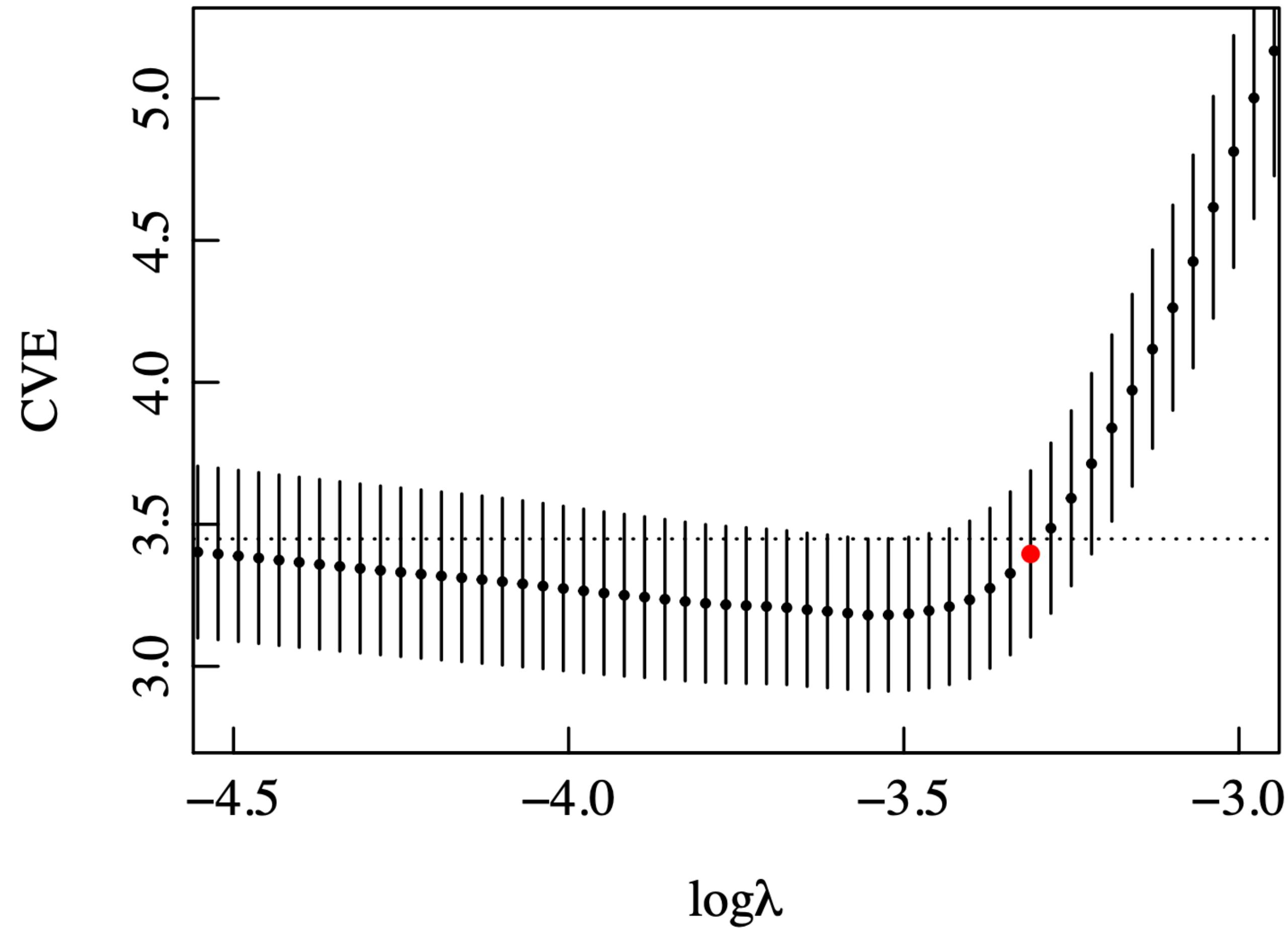
- Powerful in situations where we have more features than observations and when we believe that  $\beta$  is sparse = only a few features are truly important.
- The level of sparsity is determined by  $\lambda$ .



# The optimal model is determined by $\lambda$

## Choosing $\lambda$ correctly is crucial

- Choose  $\lambda$  that maximizes the model's generalization performance.
- Cross-validation is a straightforward and widely used method.
- Right figure: Example from a practical case discussed later
  - If  $\lambda$  is too large, the resulting model may become excessively sparse, which means it oversimplifies the data, leading to a large discrepancy between the model and the actual data.
  - if  $\lambda$  is too small, the model may overfit, capturing noise in the data, and thus reducing its ability to generalize to new, unseen data.



# Estimate the period of a variable star from time-series data

- Estimate a sparse power spectrum from data.
- The observed time-series data  $\mathbf{y}$  is modeled as a linear combination of cos and sin components at each frequency.
- A total of  $2M$  components (both sin and cos) for  $M$  different frequencies.
- Number of data points =  $N$
- According to the sampling theorem, if  $N = 2M$ , the problem is uniquely solvable.
- If some data points are missing or the sampling intervals are non-uniform → Underdetermined system with  $N < 2M$

$$\begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} = \begin{pmatrix} \cos(2\pi t_1 \nu_1) & \cdots & \cos(2\pi t_1 \nu_M) \\ \vdots & \ddots & \vdots \\ \cos(2\pi t_N \nu_1) & \cdots & \cos(2\pi t_N \nu_M) \end{pmatrix} \begin{pmatrix} a_1 \\ \vdots \\ a_M \\ b_1 \\ \vdots \\ b_M \end{pmatrix}$$

$$\mathbf{y} = \mathcal{F}^{-1} \boldsymbol{\beta}$$

# Solve using Group LASSO

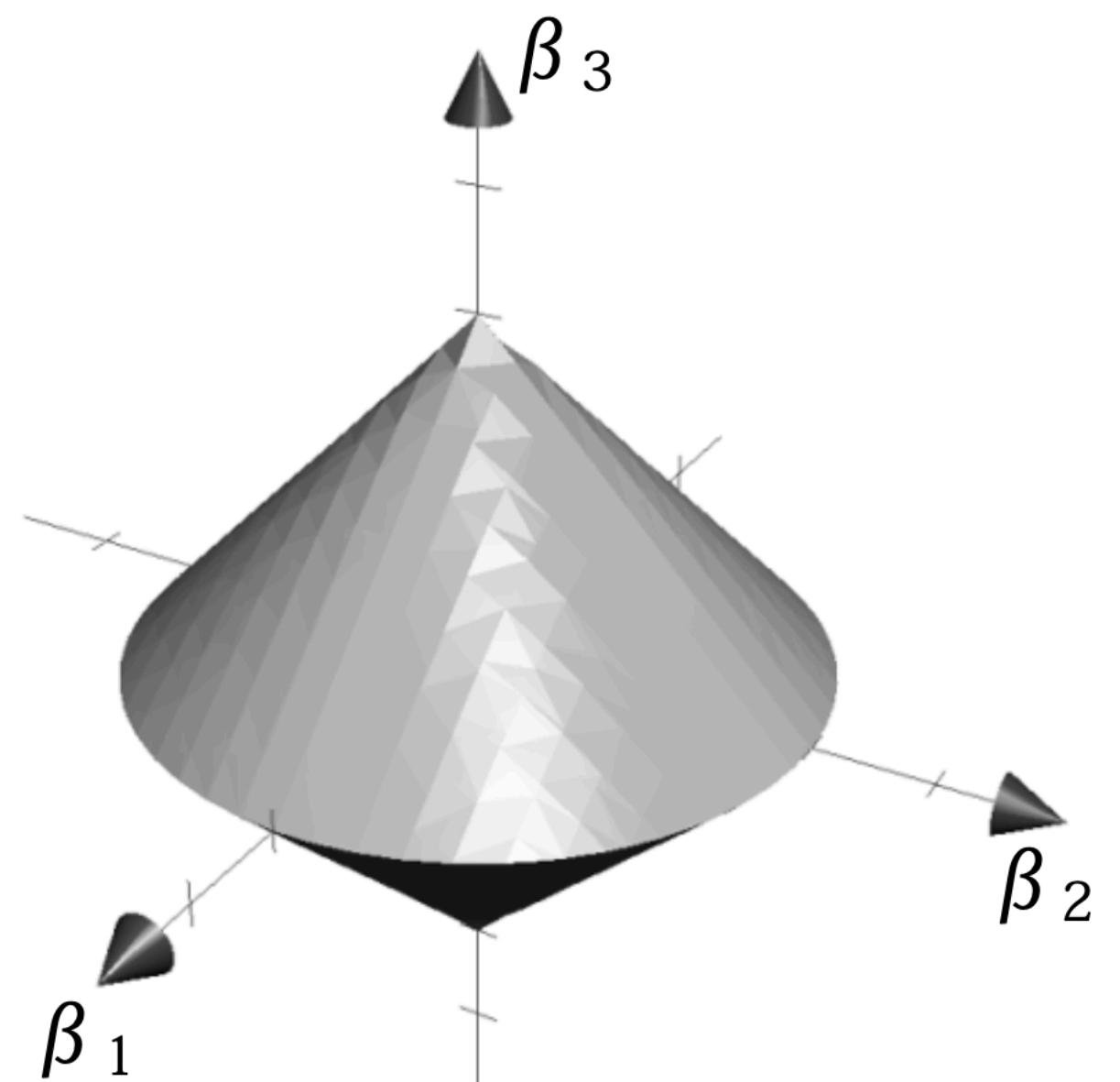
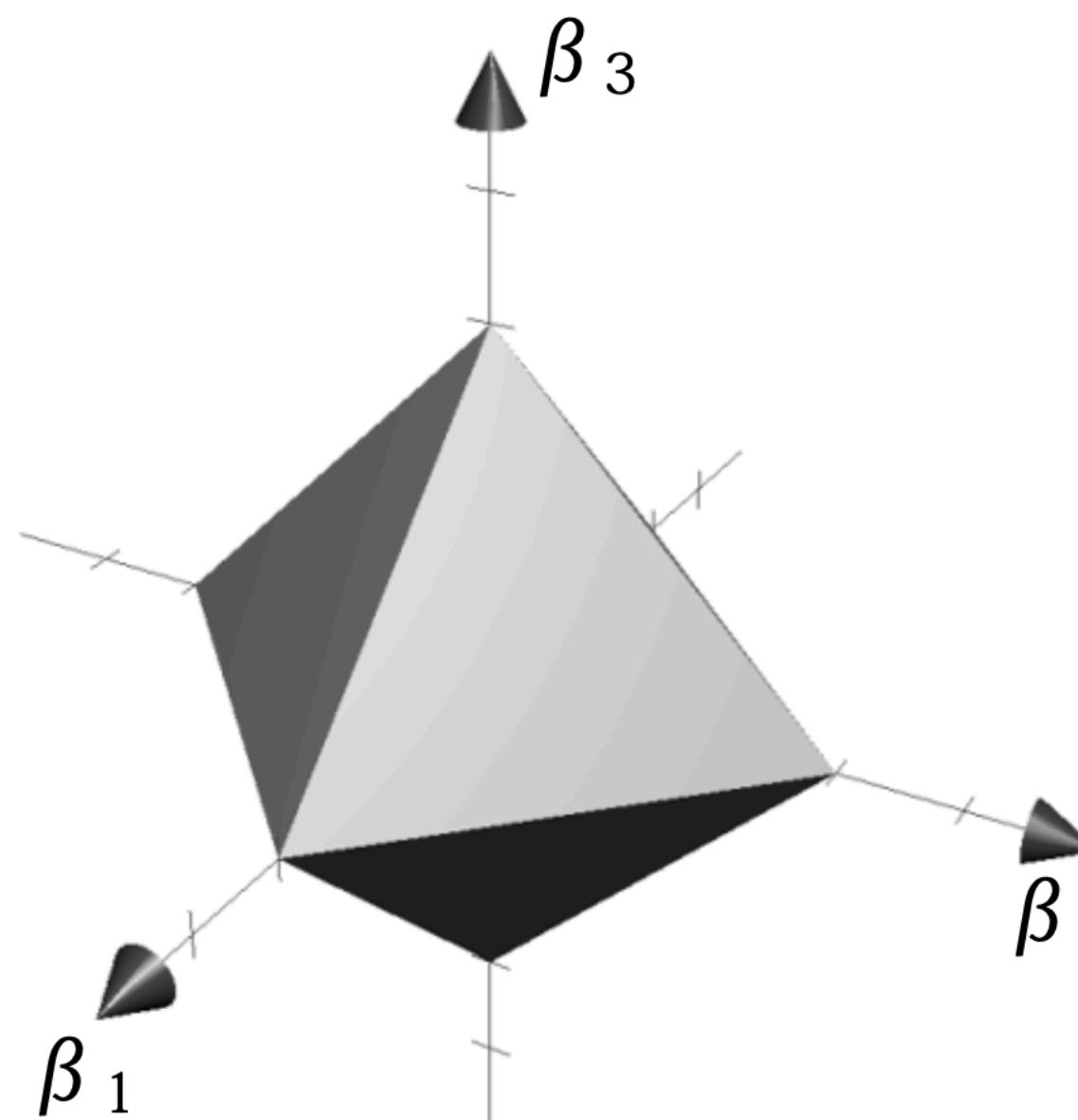
- Assume that the power spectrum is sparse.
- Using standard LASSO, the sin and cos coefficients for each frequency would be treated independently (left figure).

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathcal{F}^{-1}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1$$

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathcal{F}^{-1}\boldsymbol{\beta}\|_2^2 + \lambda \sum_i \sqrt{a_i^2 + b_i^2}$$

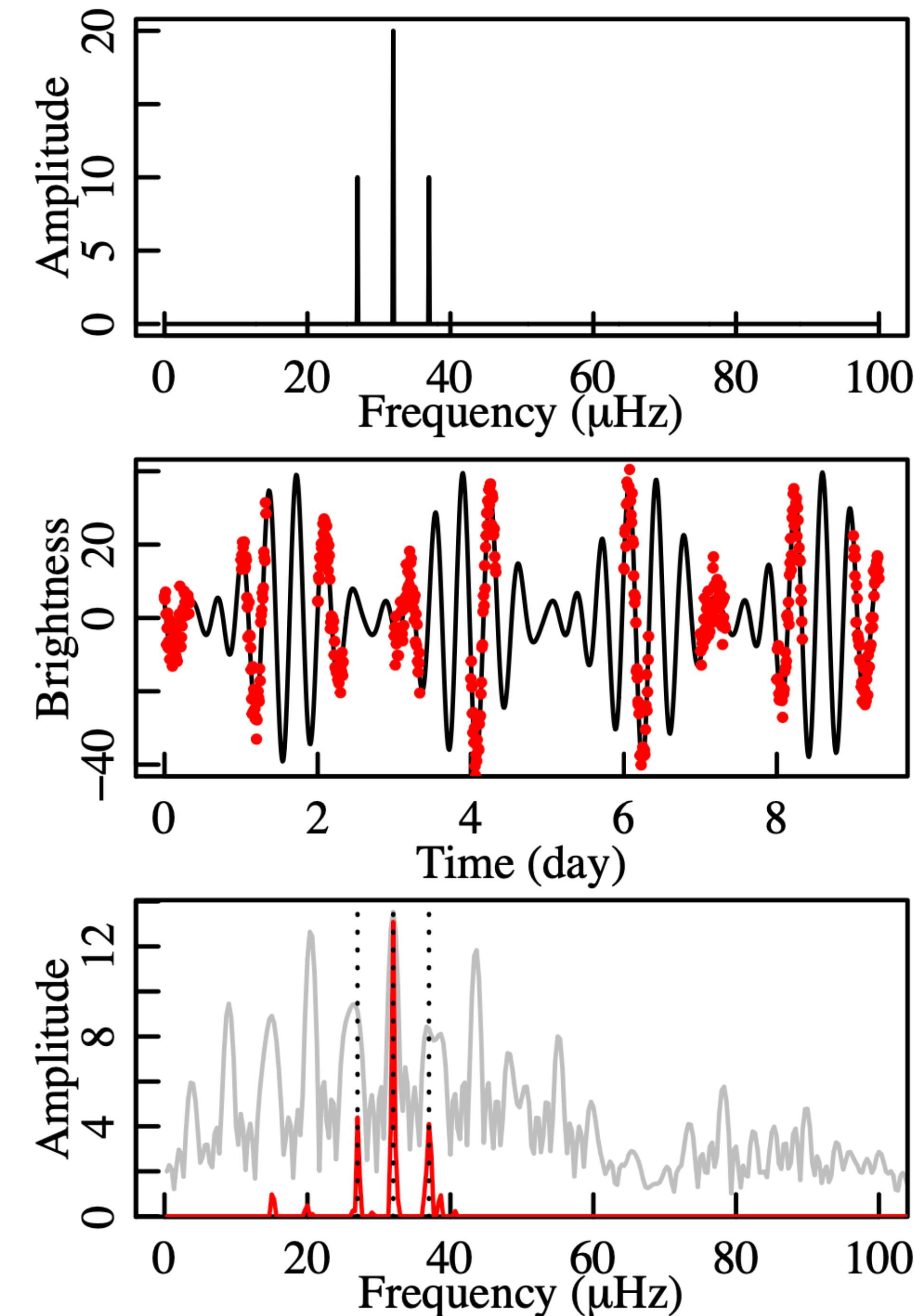
- With Group LASSO, coefficients are grouped by frequency.

- Bottom-right figure:  $\beta_1$  and  $\beta_2$ , representing the sine and cosine components for a particular frequency, are grouped together.
- Group LASSO enforces sparsity at the group level, meaning that either the entire group (both  $\beta_1$  and  $\beta_2$ ) or another coefficient, such as  $\beta_3$ , is set to zero.



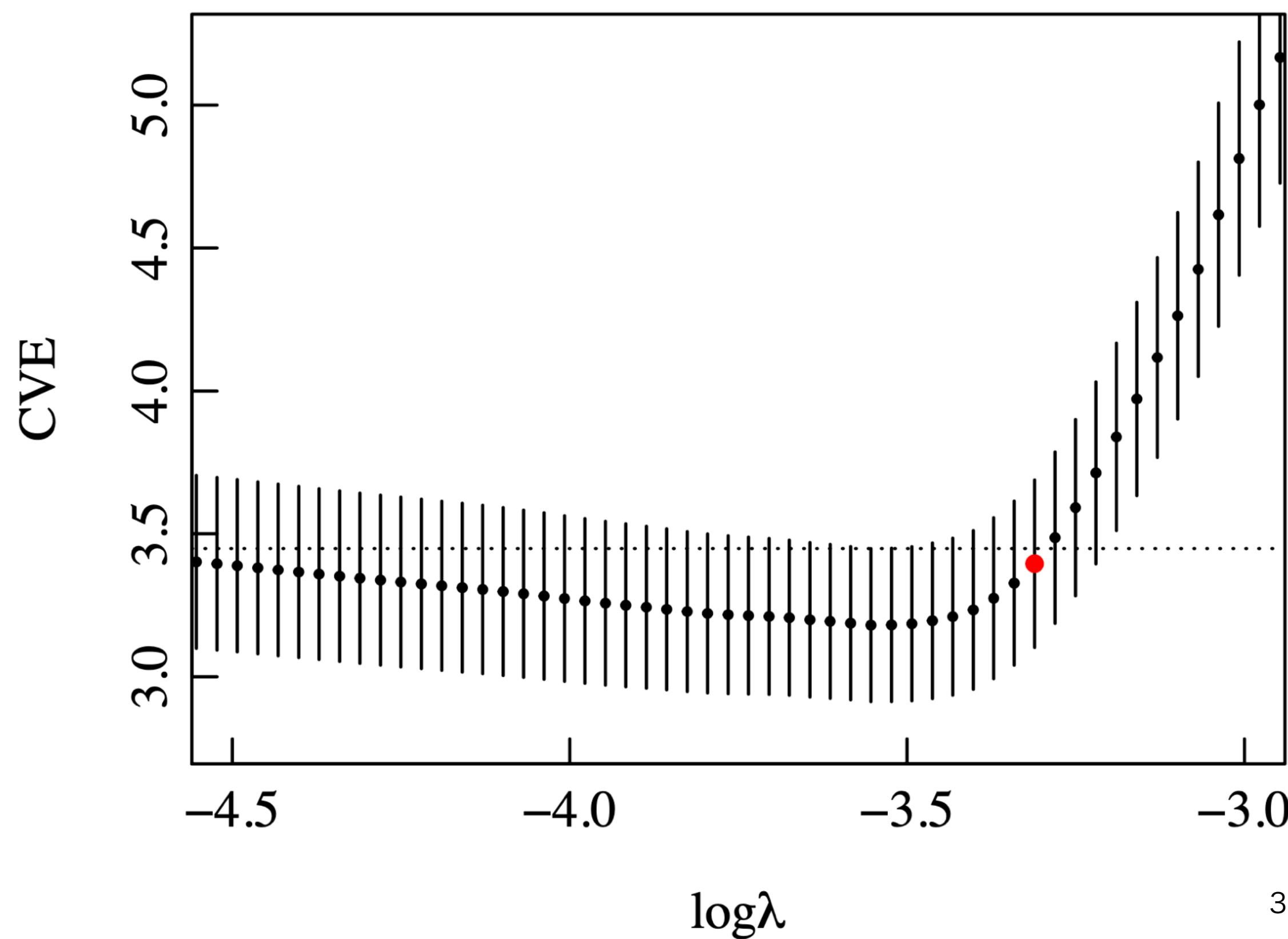
# Experiments

- Assume three distinct frequency signals (top panel)
- Perform an inverse Fourier transform to create a time series dataset (black line in the middle panel)
- Downsample the data and add Gaussian noise (red dots)
- Perform a standard Fourier transform with the red dots → the gray power spectrum (bottom panel)
  - While the strongest signal is detected, the other weaker signals are buried under aliasing
  - Perform a group LASSO regression → the red curve
    - Aliasing is significantly reduced, and the true signals become much clearer.

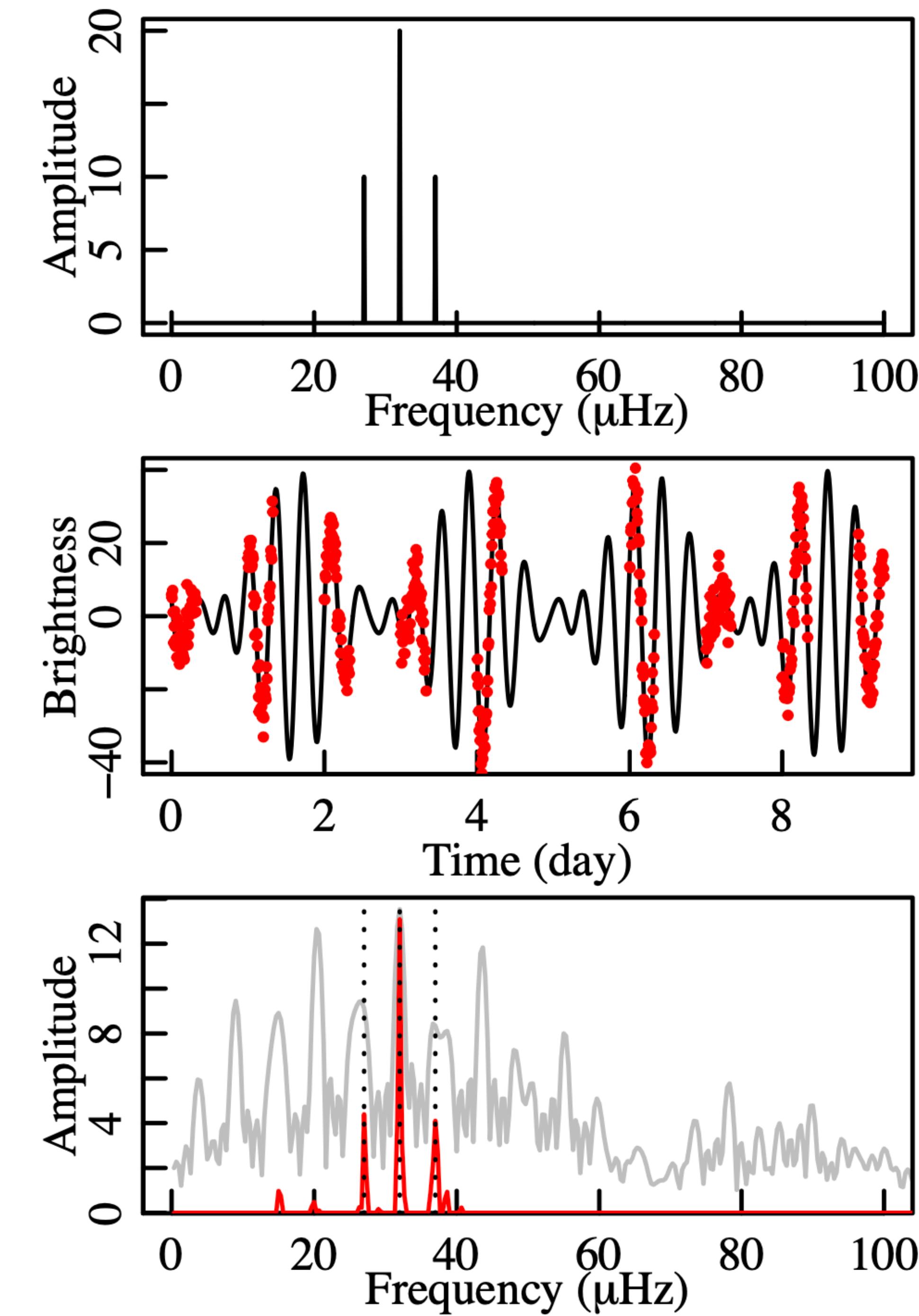


# Experiments

- The regularization parameter  $\lambda$  is determined using cross-validation and the **one-standard-error rule**

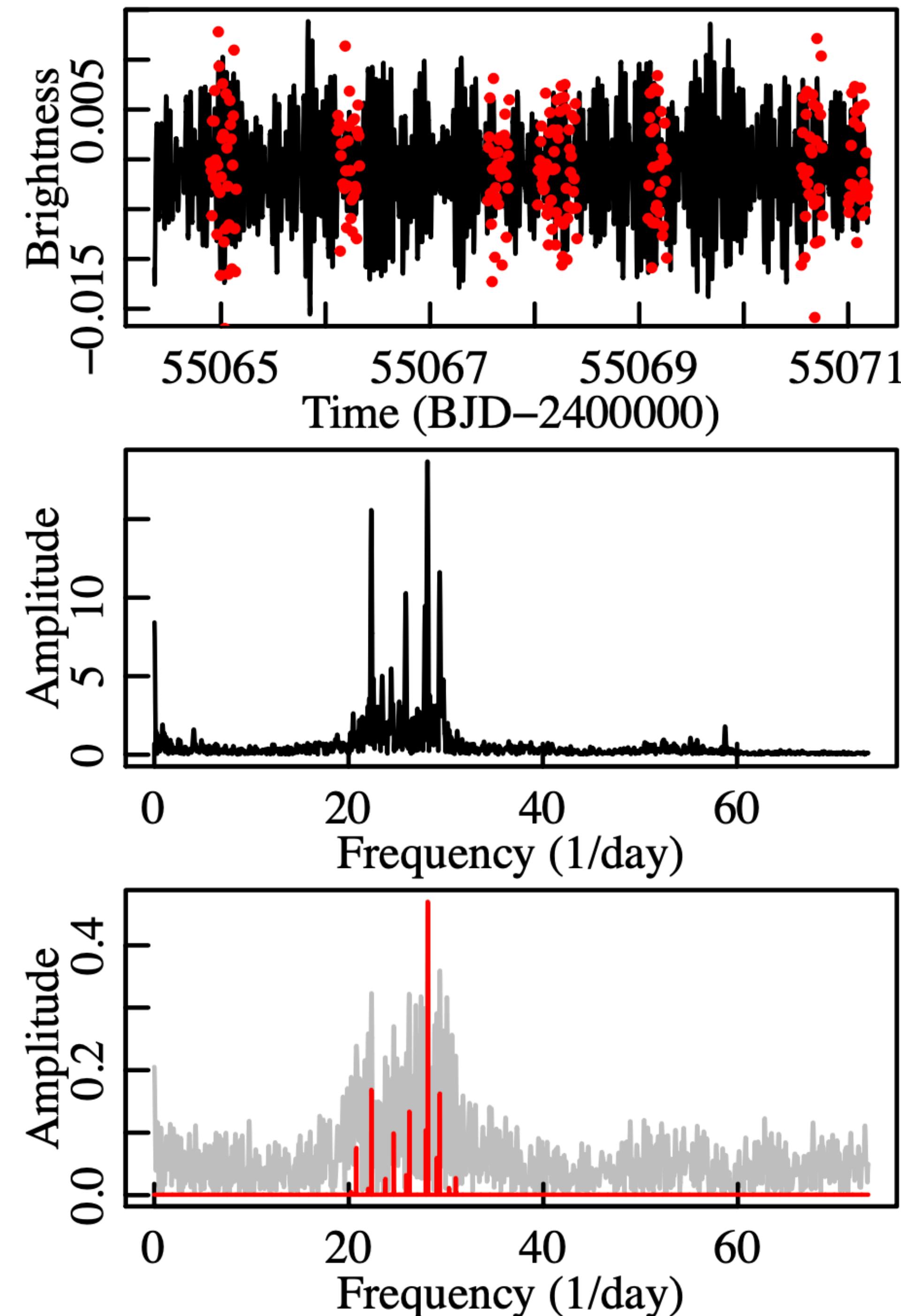


31

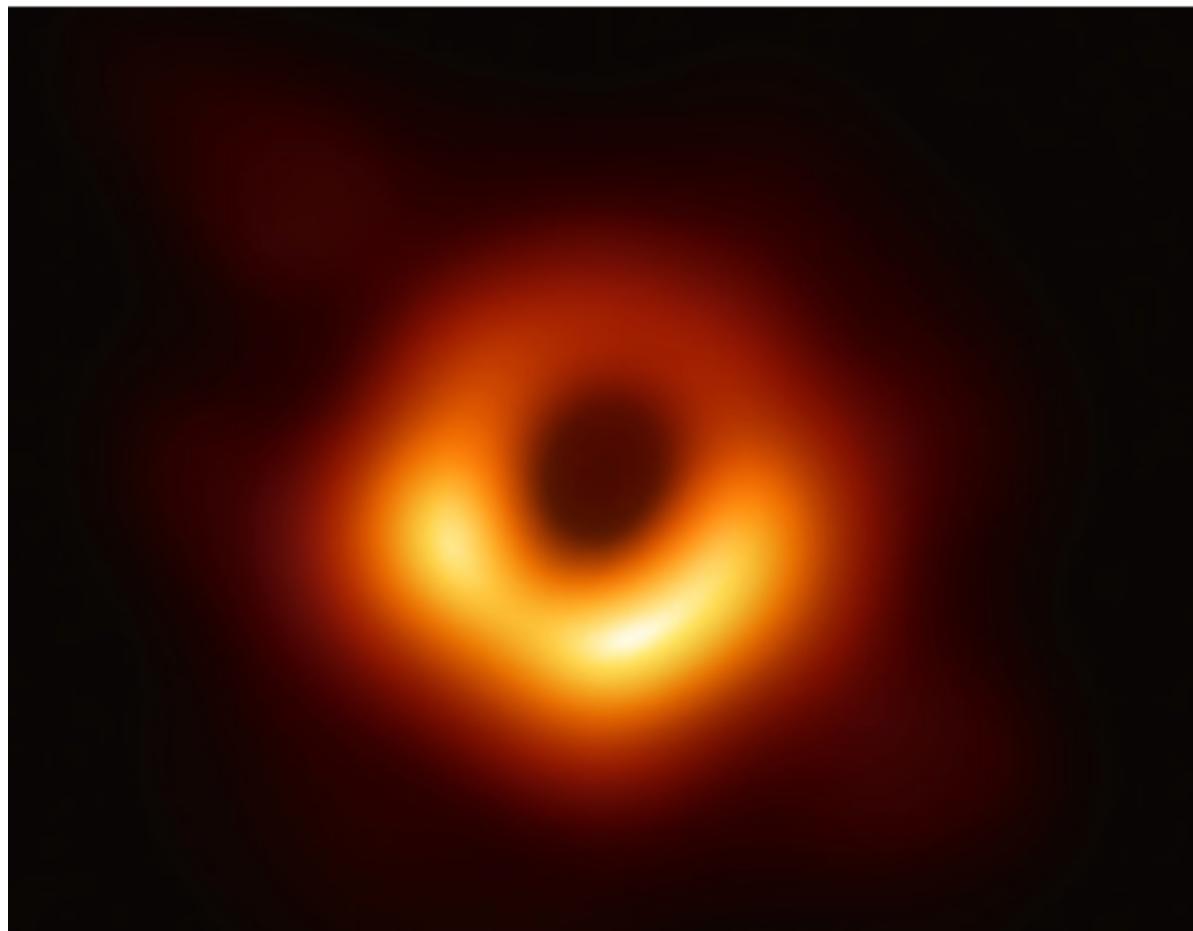


# Apply to real data

- $\delta$  Sct type pulsating star observed with the Kepler satellite (black line in the top panel)
  - Nearly complete data. Obtaining such comprehensive data from ground-based observations is challenging.
  - The power spectrum is nearly complete (middle panel)
- Simulating ground-based observations, downsample the data and add noise (red dots in the top panel)
- Performing a standard Fourier transform results in the gray line in the bottom panel → Full of aliasing artifacts
- Results from Group LASSO (red line in the bottom panel)
  - Group LASSO effectively detects the strong signals present in the original power spectrum.



# Event Horizon Telescope, again as a regularized least-squares



©EHT

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \left[ \frac{1}{2} \|\mathbf{v} - F\mathbf{x}\|_{\ell_2}^2 + \lambda_1 \|\mathbf{x}\|_{\ell_1} + \lambda_{\text{TSV}} \text{TSV}(\mathbf{x}) \right],$$

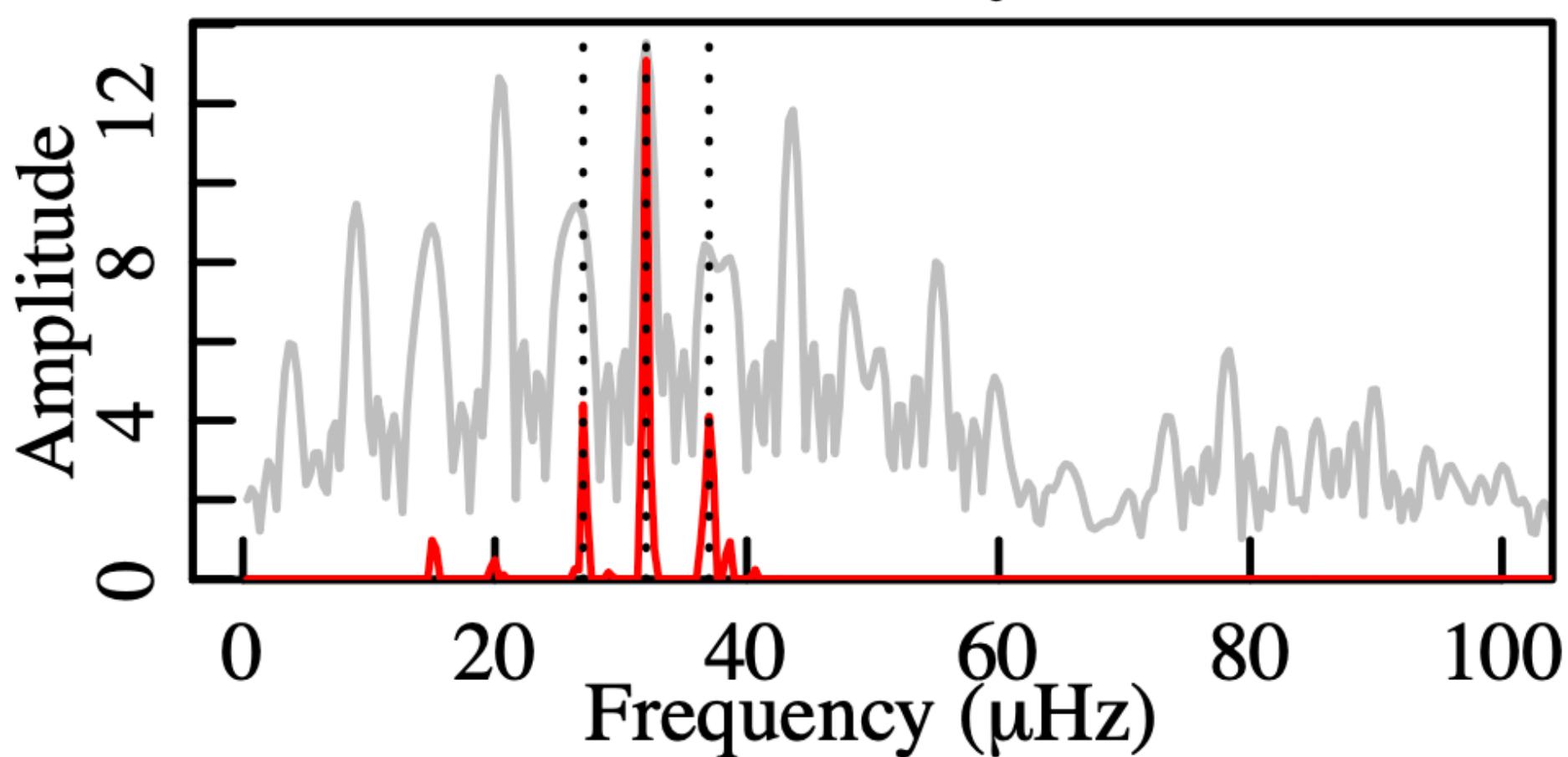
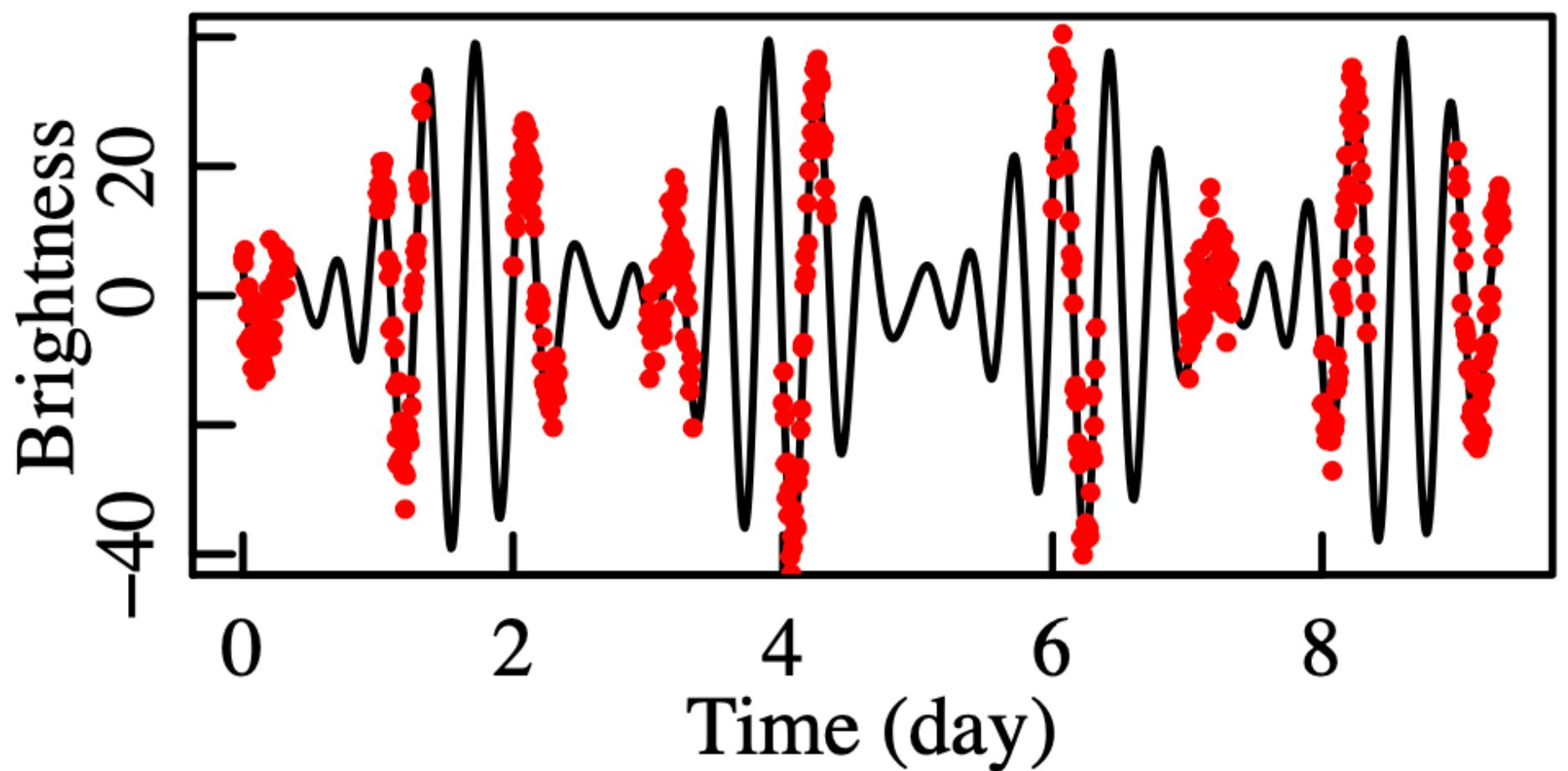
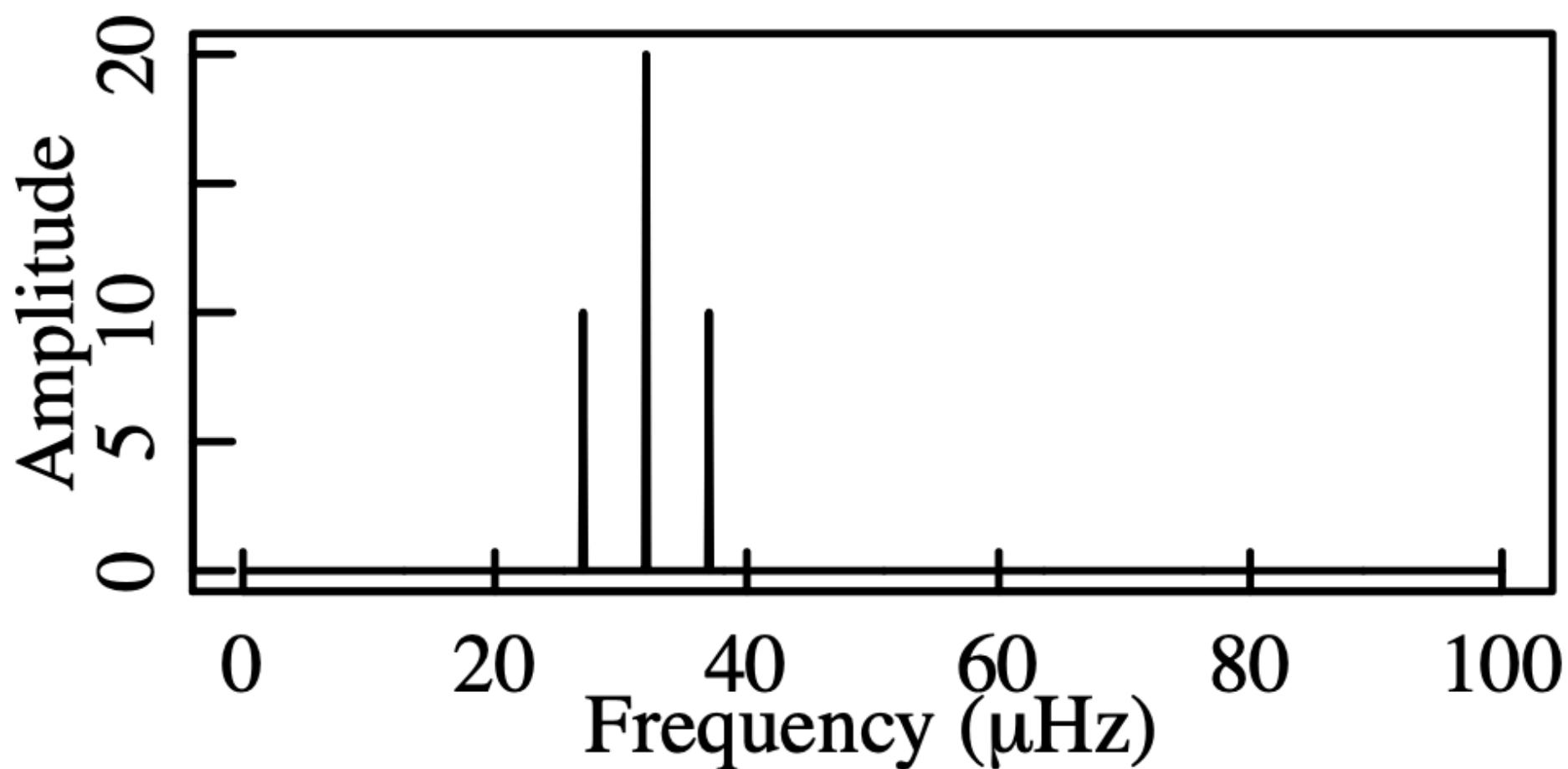
$$\text{TSV}(\mathbf{x}) = \sum_{ij} \left[ (x_{ij} - x_{i,j-1})^2 + (x_{ij} - x_{i-1,j})^2 \right].$$

Event Horizon Telescope Collaboration (2019)

# Hands-on exercise #2

Power spectrum estimation with LASSO

Lecture\_Day1\_Uemura/02\_LASSO.ipynb

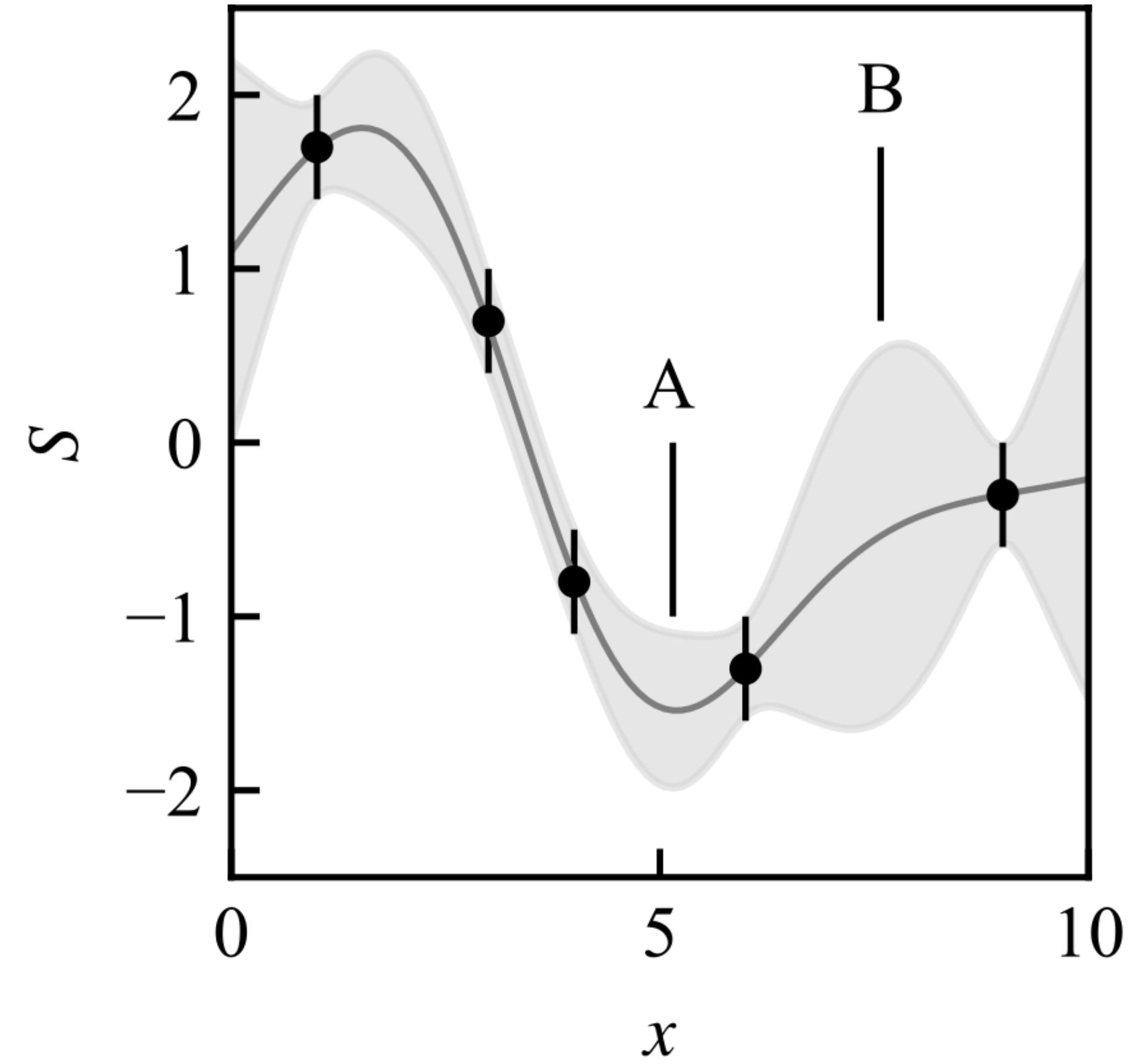


# 3. Gaussian processes

**Key words:** Kernel function, Bayesian optimization

# What's Gaussian processes (GP) useful for

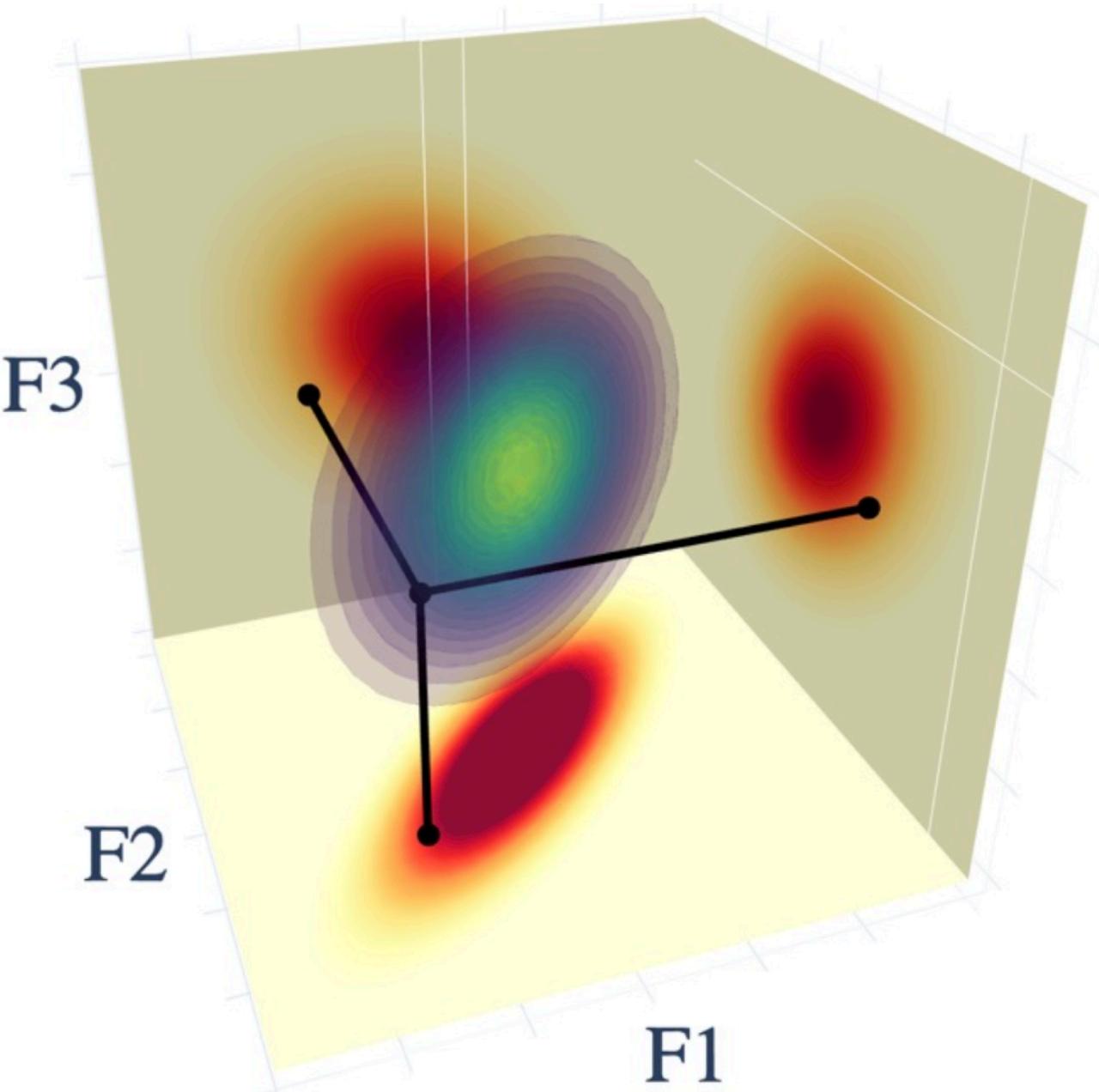
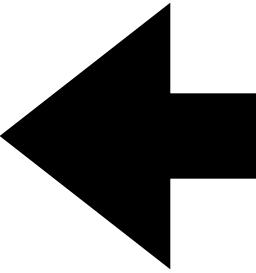
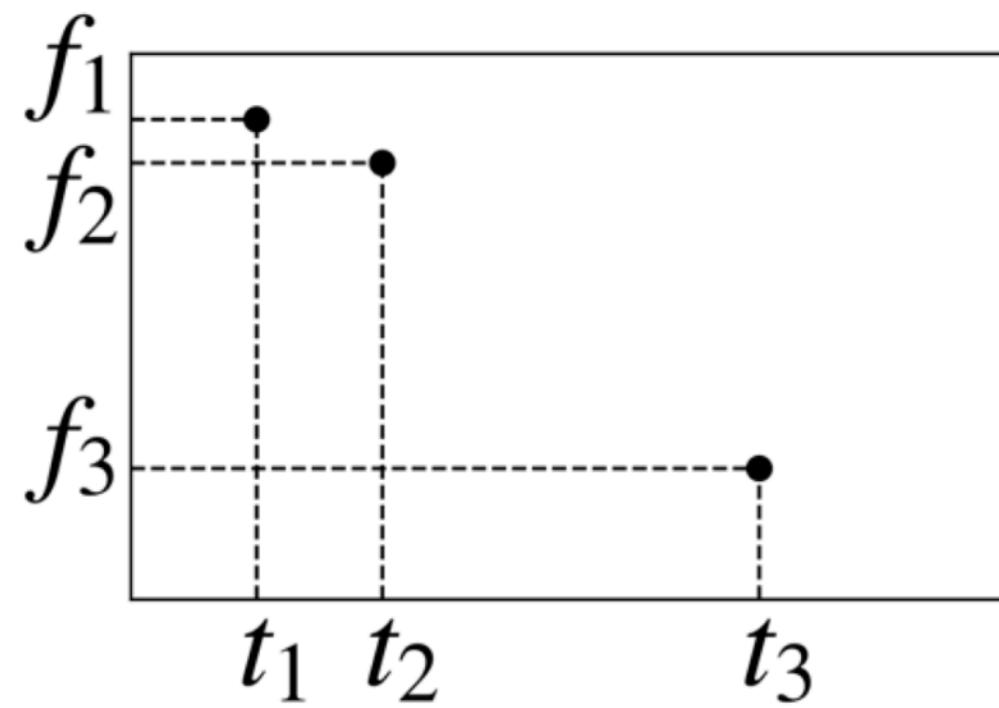
- Providing plausible curves (or surfaces) when the functional form of the data is unknown
  - Useful for interpolation tasks
  - Also offering plausible uncertainty regions
    - Useful for finding the maxima and minima of functions and constructing surrogate models
  - Computationally efficient
    - Relatively straightforward to implement from scratch
    - Python provides convenient modules for working with GPs, such as GPy and sklearn.gaussian\_process.



Example of Gaussian Process regression. If you want to find the minimum value of the function, you should measure at point "A" next. If you want to create a more accurate surrogate model, you should measure at point "B" next.

# Definition of Gaussian Processes

Consider the N data points as a sample from an N-dimensional normal distribution



- Set the mean  $\mu$  to zero by data preprocessing
- You need to define the covariance matrix = a core of GPs.
- Example: When there are 3 data points

$$\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \sigma_{13} \\ \sigma_{21} & \sigma_{22} & \sigma_{23} \\ \sigma_{31} & \sigma_{32} & \sigma_{33} \end{pmatrix}$$

$$p(\mathbf{z}) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp \left\{ -\frac{1}{2} (\mathbf{z} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{z} - \boldsymbol{\mu}) \right\}$$

- Each element is a function of  $\sigma_{12} = k(t_1, t_2)$
- If there are N data points, you need  $N(N-1)/2$  elements → Too many parameters

$$p(\mathbf{z}) = \mathcal{N}(\boldsymbol{\mu}, \Sigma)$$

# Kernel functions

- Replacing covariances to the kernel functions:  $\sigma_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$
- RBF (Gaussian) kernel
  - . The distance between two points  $\mathbf{x}_i, \mathbf{x}_j \sim$  their correlation.
  - . Powerful for capturing smooth and continuous relationships in the data.
  - .  $\theta_2$  = length-scale
- Estimating a few kernel parameters from N data points → Gaussian Process Regression

$$\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \sigma_{13} \\ \sigma_{21} & \sigma_{22} & \sigma_{23} \\ \sigma_{31} & \sigma_{32} & \sigma_{33} \end{pmatrix}$$

RBF kernel

$$k(\mathbf{x}_i, \mathbf{x}_j) = \theta_1 \exp \left\{ -\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{\theta_2} \right\}$$

Exponential kernel  $k(\mathbf{x}_i, \mathbf{x}_j) = \theta_1 \exp \left\{ -\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2}{\theta_2} \right\}$

Periodic kernel

$$k(\mathbf{x}_i, \mathbf{x}_j) = \theta_1 \exp \left\{ -\theta_2 \cos \frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2}{\theta_3} \right\}$$

Linear kernel

$$k(\mathbf{x}_i, \mathbf{x}_j) = \theta_1 \mathbf{x}_i^T \mathbf{x}_j$$

Matern kernel

$$K(x, x') = \frac{2^{1-\nu}}{\Gamma(\nu)} \left( \frac{\sqrt{2\nu}|x-x'|}{\rho} \right)^\nu K_\nu \left( \frac{\sqrt{2\nu}|x-x'|}{\rho} \right)$$

$K_\nu$  is the modified Bessel function of the second kind

# Gaussian Process Regression

## Estimation of the kernel parameters from data

Likelihood

$$p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{f}, \sigma^2 \mathbf{I})$$

Prior

$$p(\mathbf{f}) = \mathcal{N}(\mathbf{0}, \mathbf{K})$$

- A Bayesian model where a sample  $\mathbf{f}$  from a Gaussian process defined by the covariance matrix  $\mathbf{K}$  is subject to measurement error  $\sigma$ .

Marginalized likelihood

$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f})d\mathbf{f}$$

$$= \mathcal{N}(\mathbf{0}, \mathbf{C}) = \frac{1}{\sqrt{(2\pi)^N |\mathbf{C}|}} \exp \left\{ -\frac{1}{2} \mathbf{y}^T \mathbf{C}^{-1} \mathbf{y} \right\}$$

$$\mathbf{C} = \mathbf{K} + \sigma^2 \mathbf{I}$$

- Marginalize over  $\mathbf{f}$ .
- Even with Gaussian noise added to the Gaussian process, it remains a Gaussian process.

Log-likelihood

$$\log p(\mathbf{y}) = -\frac{N}{2} \log(2\pi) - \frac{1}{2} \log |\mathbf{C}| - \frac{1}{2} \mathbf{y}^T \mathbf{C}^{-1} \mathbf{y}$$

- Maximize the log (marginal) likelihood using gradient methods or similar → Determines the kernel parameters.

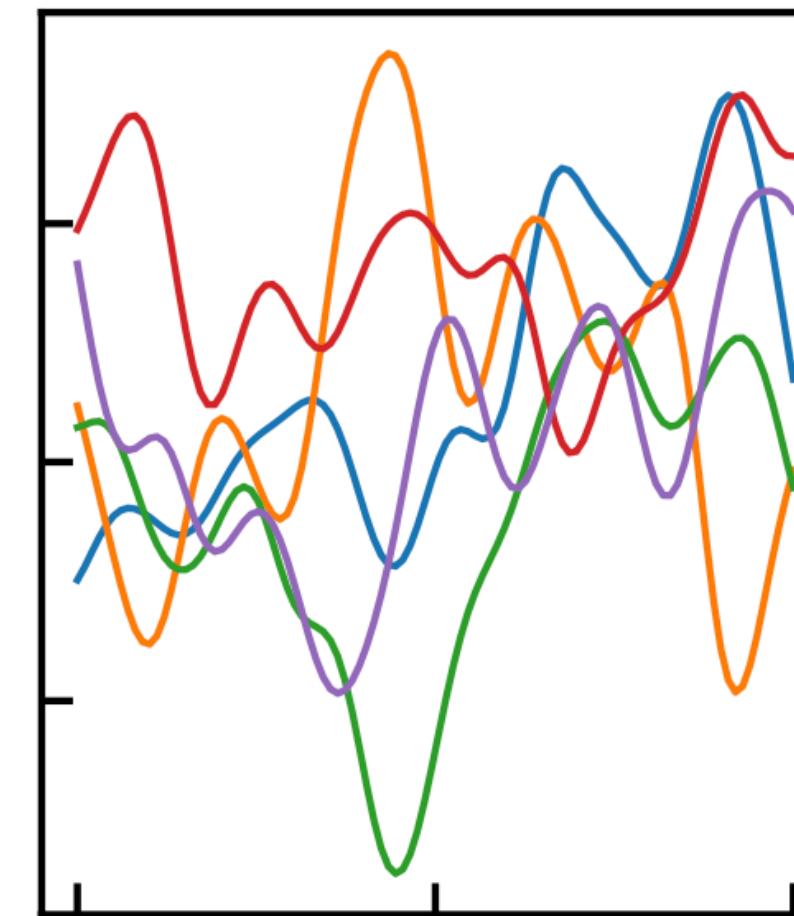
# Prior and Posterior of kernels

$$p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{f}, \sigma^2 \mathbf{I})$$

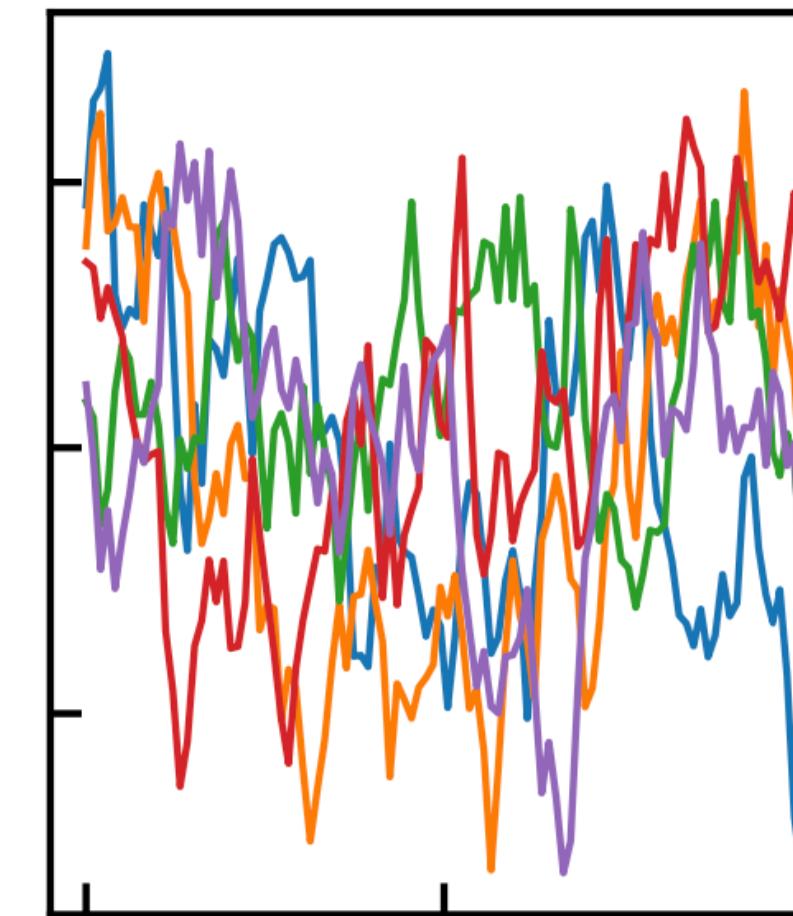
$$p(\mathbf{f}) = \mathcal{N}(\mathbf{0}, \mathbf{K})$$

- **Top:** Samples from a Gaussian process with a given set of kernel parameters. = Prior  $p(\mathbf{f})$
- **Bottom:** Predictions under the conditions where data has been obtained. = Posterior  $p(\mathbf{f}|\mathbf{y})$

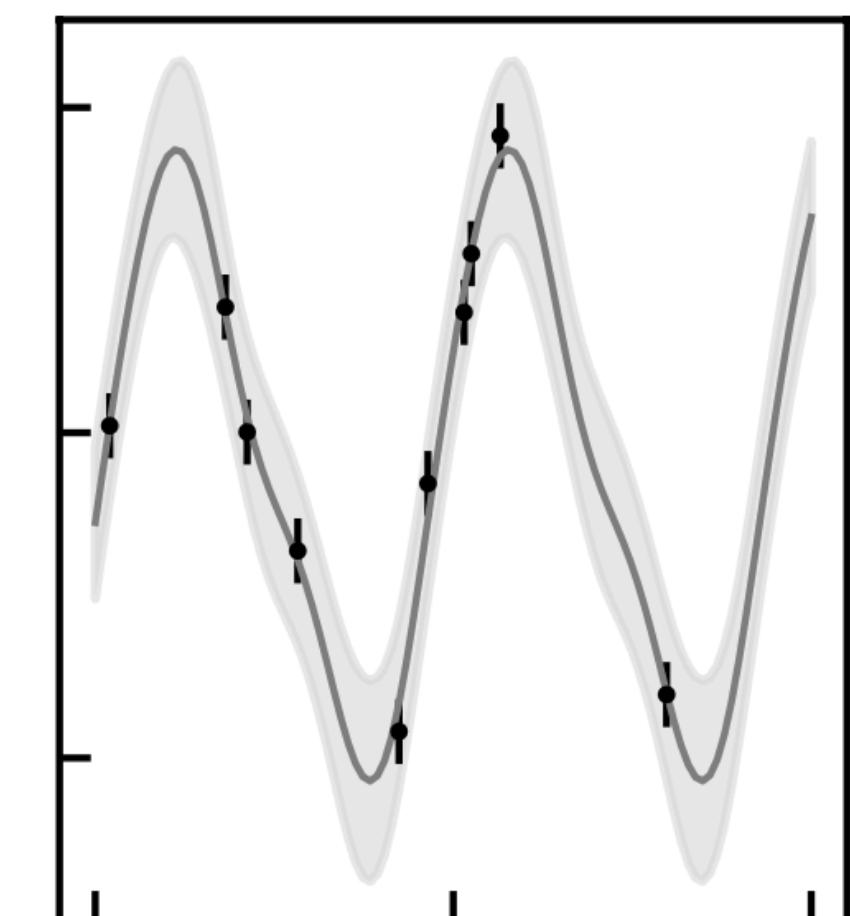
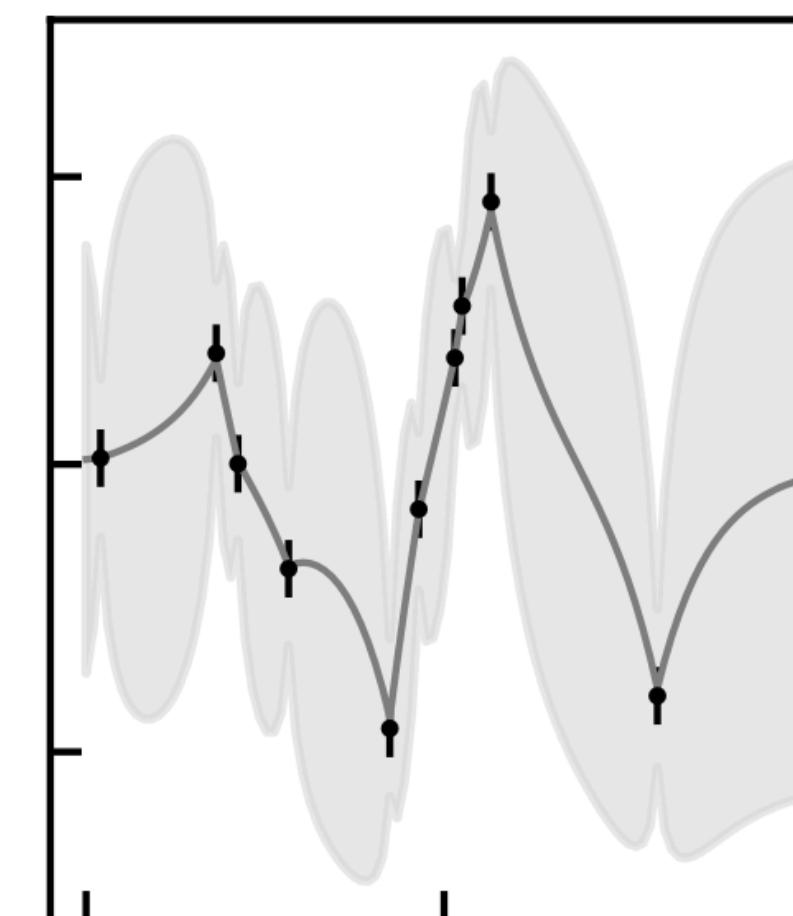
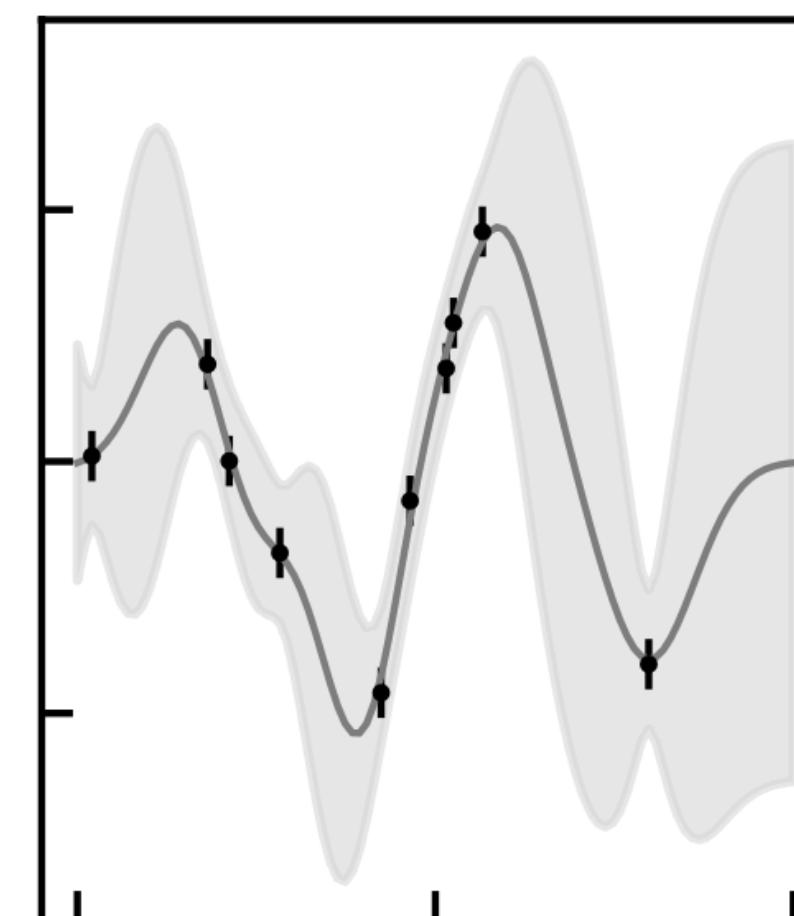
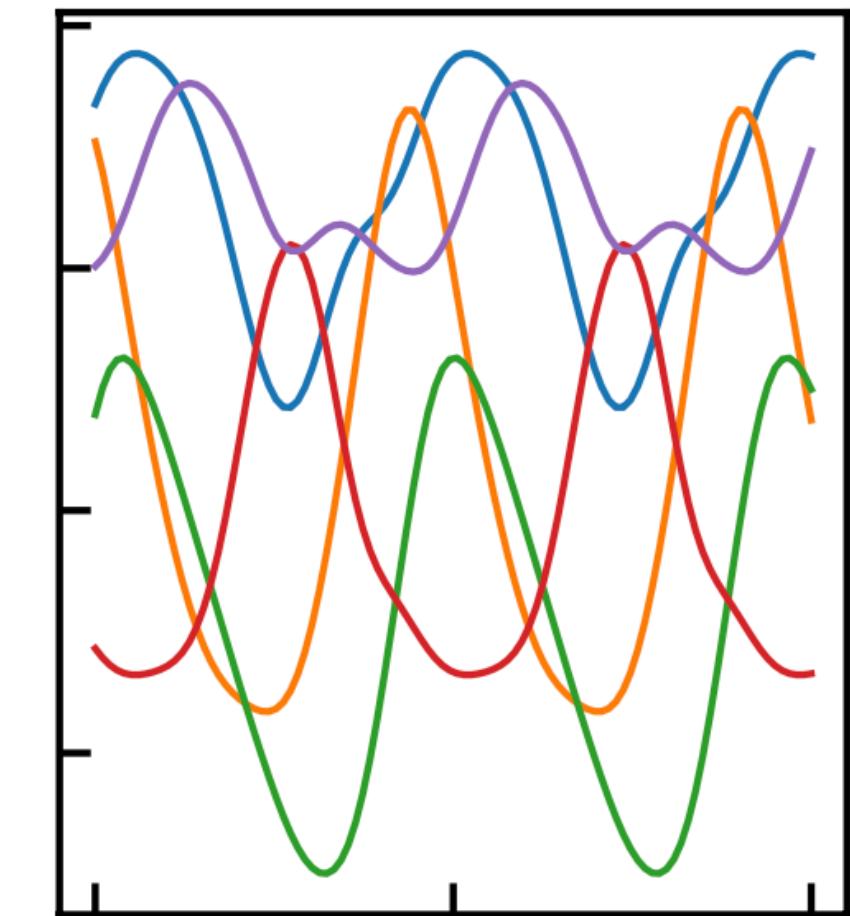
RBF kernel



Exp. kernel



Periodic kernel



# Prediction of $y^*$ at $x^*$

Joint Prob.  $p(\mathbf{y}' = (\mathbf{y}, y_*))$

Same GP  $p(\mathbf{y}') = \mathcal{N}(\mathbf{0}, \mathbf{C}')$

$$\mathbf{C}' = \begin{pmatrix} \mathbf{C} & \mathbf{k}_* \\ \mathbf{k}_*^T & c \end{pmatrix} \quad \begin{aligned} \mathbf{k}_*^T &= (k(\mathbf{x}_*, \mathbf{x}_1), k(\mathbf{x}_*, \mathbf{x}_2), \dots, k(\mathbf{x}_*, \mathbf{x}_N)) \\ c &= k(\mathbf{x}_*, \mathbf{x}_*) \end{aligned}$$

Theorem for  
the conditional  
multivariate  
Gaussian

$$\begin{pmatrix} \mathbf{y}^{(1)} \\ \mathbf{y}^{(2)} \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \boldsymbol{\mu}^{(1)} \\ \boldsymbol{\mu}^{(2)} \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix} \right)$$

$$p(\mathbf{y}^{(2)} | \mathbf{y}^{(1)}) = \mathcal{N}(\boldsymbol{\mu}^{(2)} + \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} (\mathbf{y}^{(1)} - \boldsymbol{\mu}^{(1)}), \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12})$$

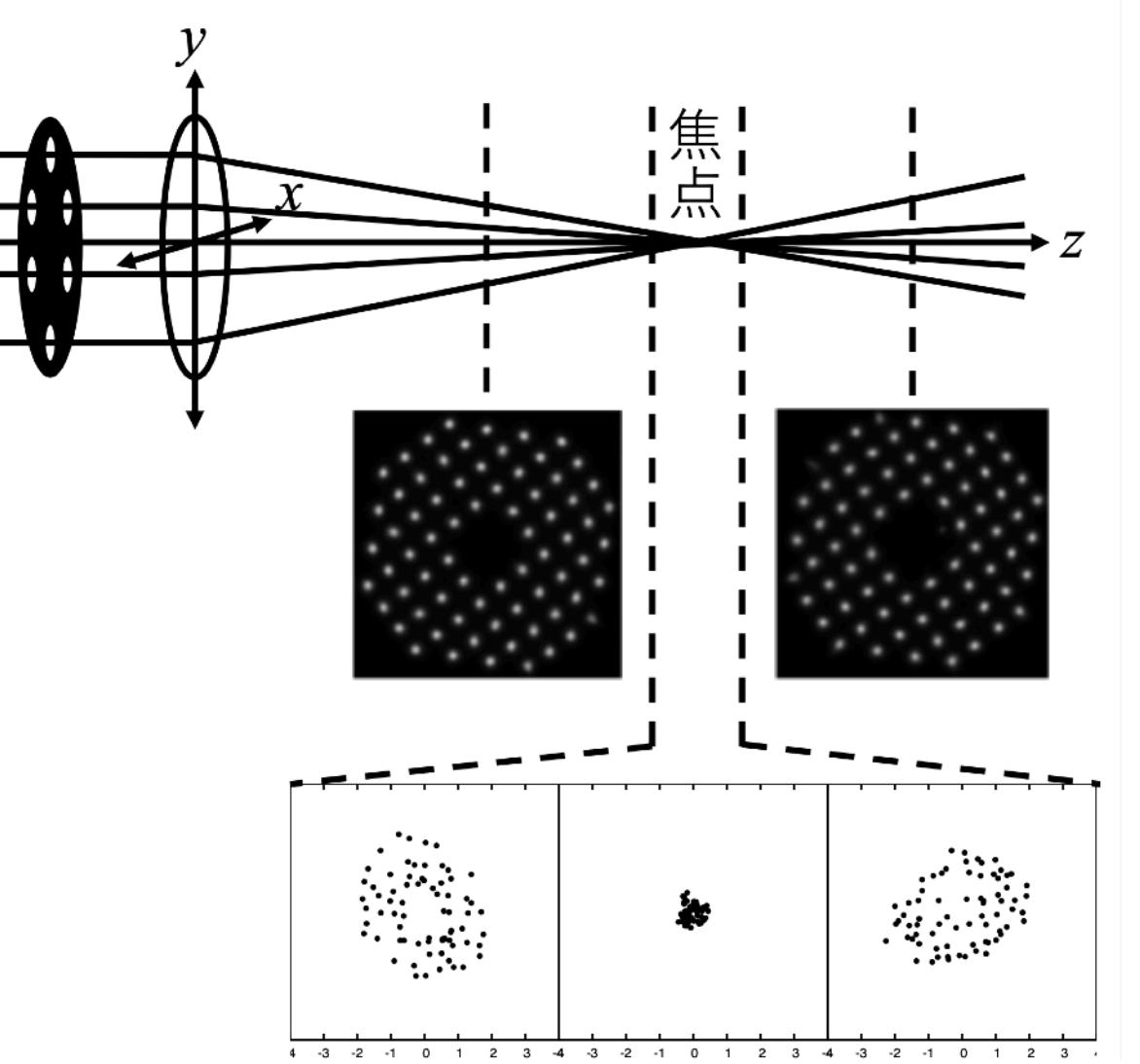
Prediction

$$p(y_* | \mathbf{y}) = \mathcal{N}(\mathbf{k}_*^T \mathbf{C}^{-1} \mathbf{y}, c - \mathbf{k}_*^T \mathbf{C}^{-1} \mathbf{k}_*)$$

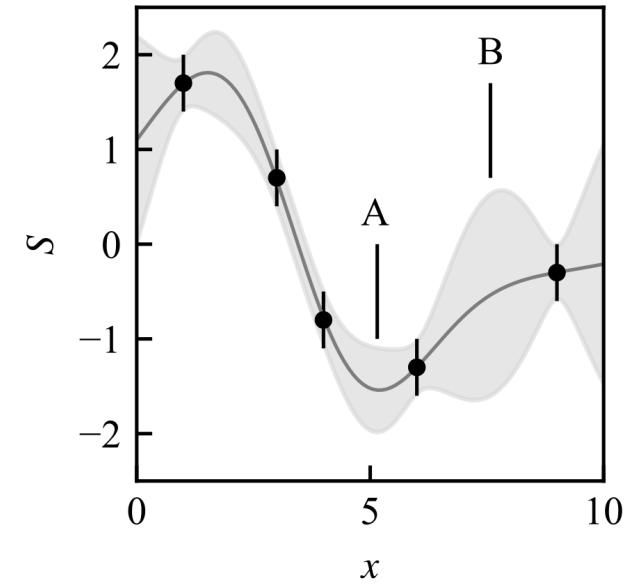
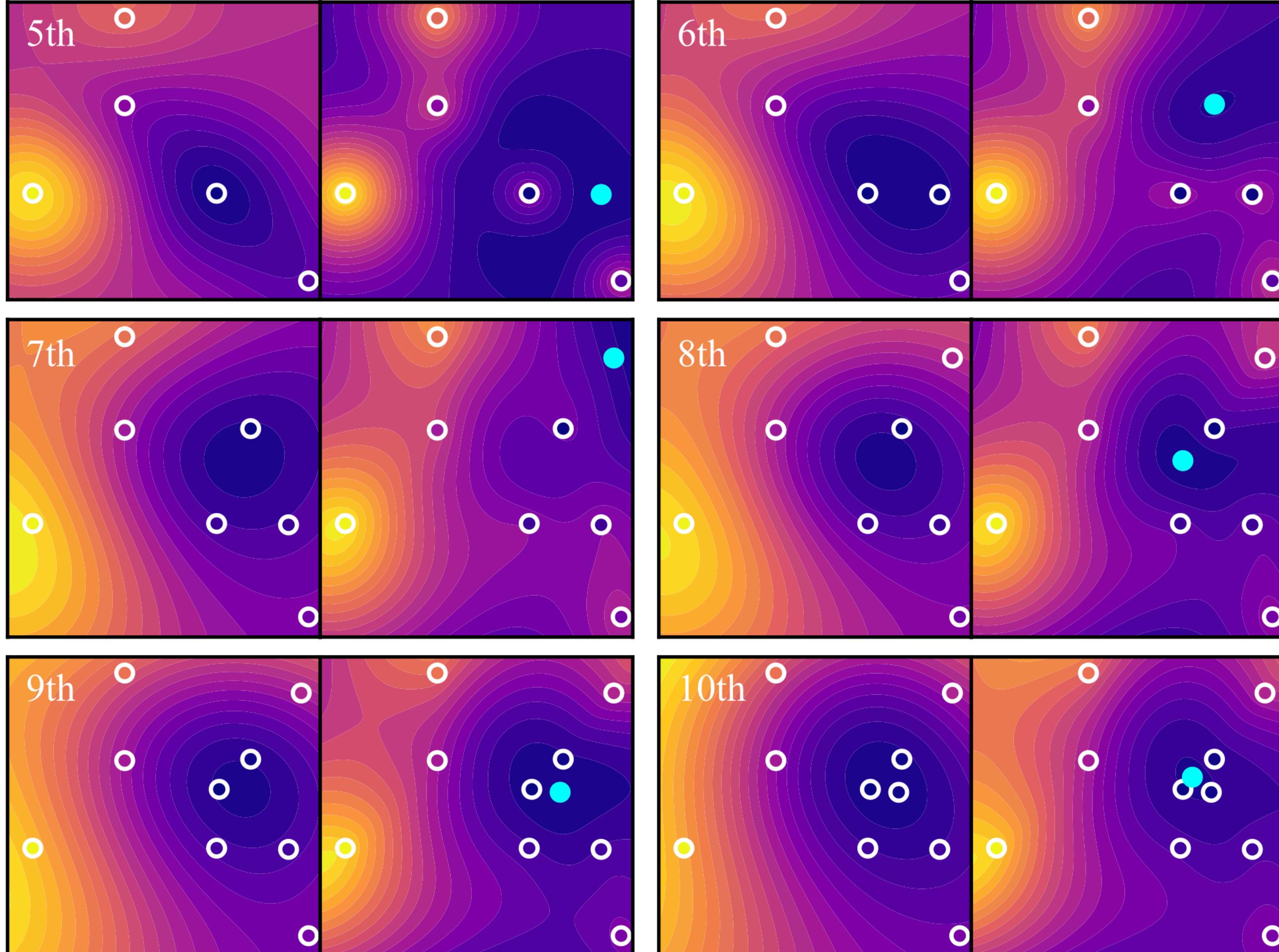
- We want to know  $p(y_* | \mathbf{y})$
- Joint probability follows the same Gaussian process.
- Use a theorem for the conditional multivariate normal distribution.
- The predictive distribution of  $y^*$  is also a normal distribution. The mean and variance can be obtained through simple matrix calculations.

# Practical Example 1

## Hartmann test for the Kanata telescope

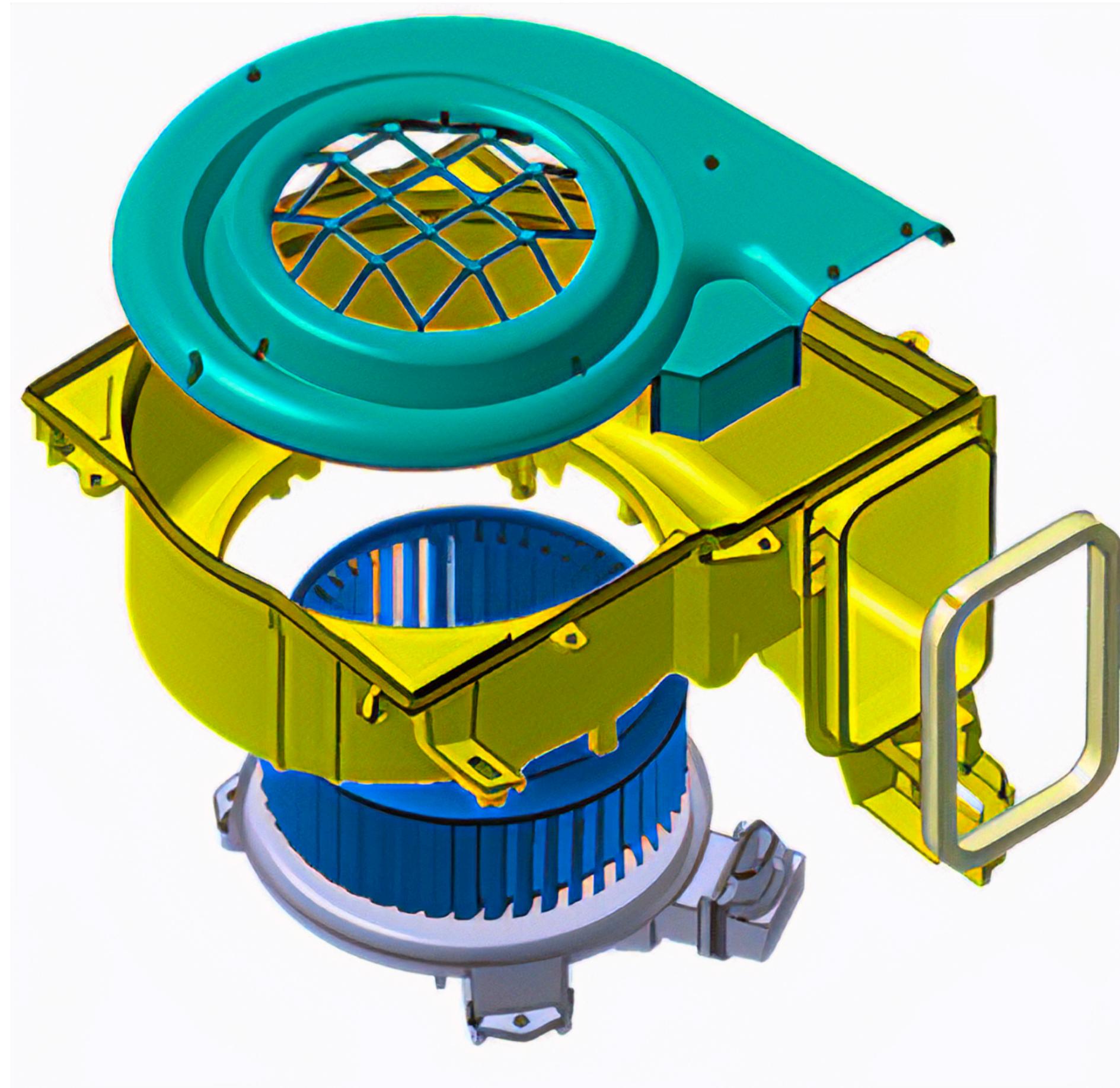


- Find the optimal  $(x, y)$  position of the secondary mirror where the light converges most effectively.
- Exhaustive search in the 2D space of  $(x, y)$  is time-consuming.
- Use “Bayesian optimization” with a Gaussian process to perform the search.



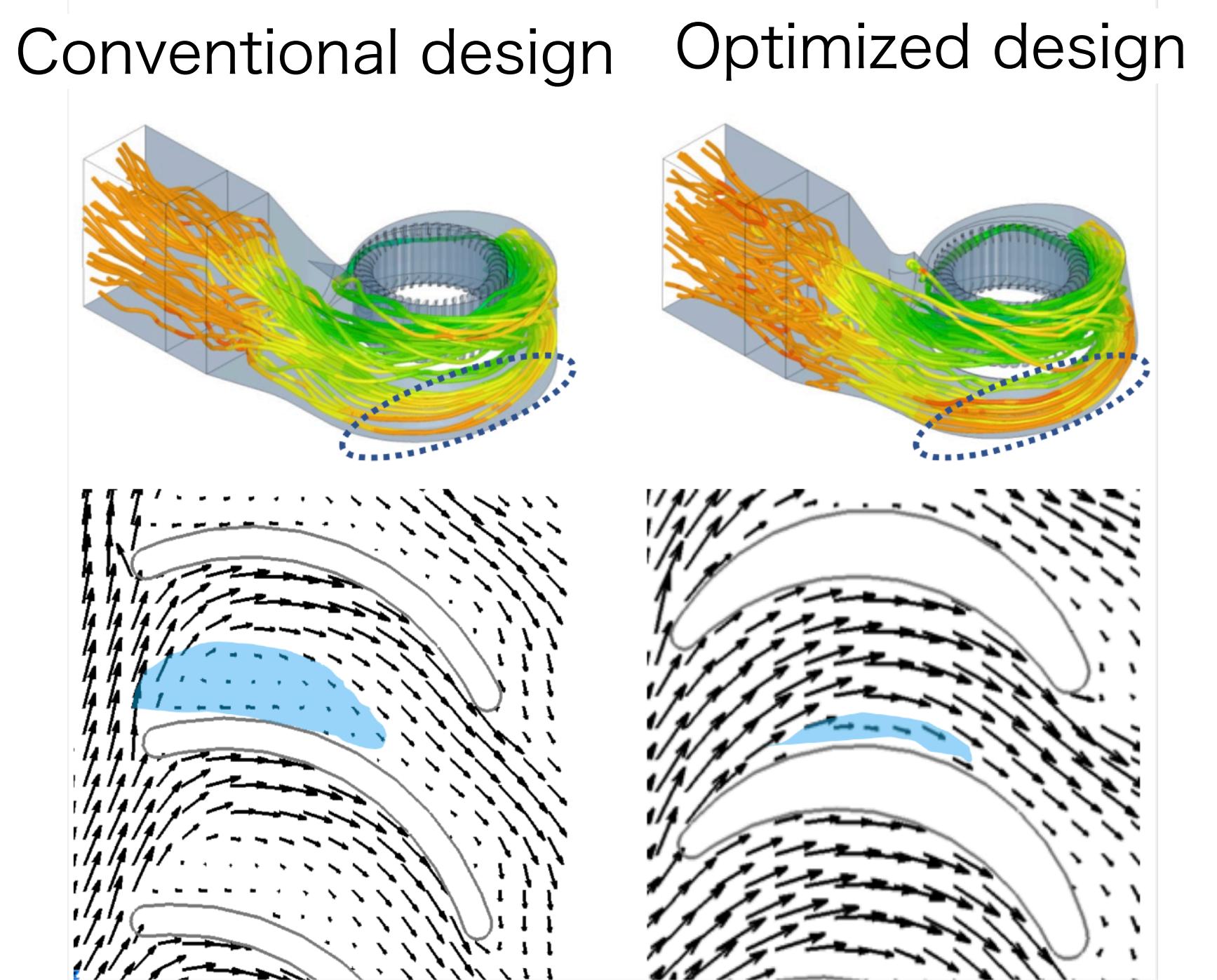
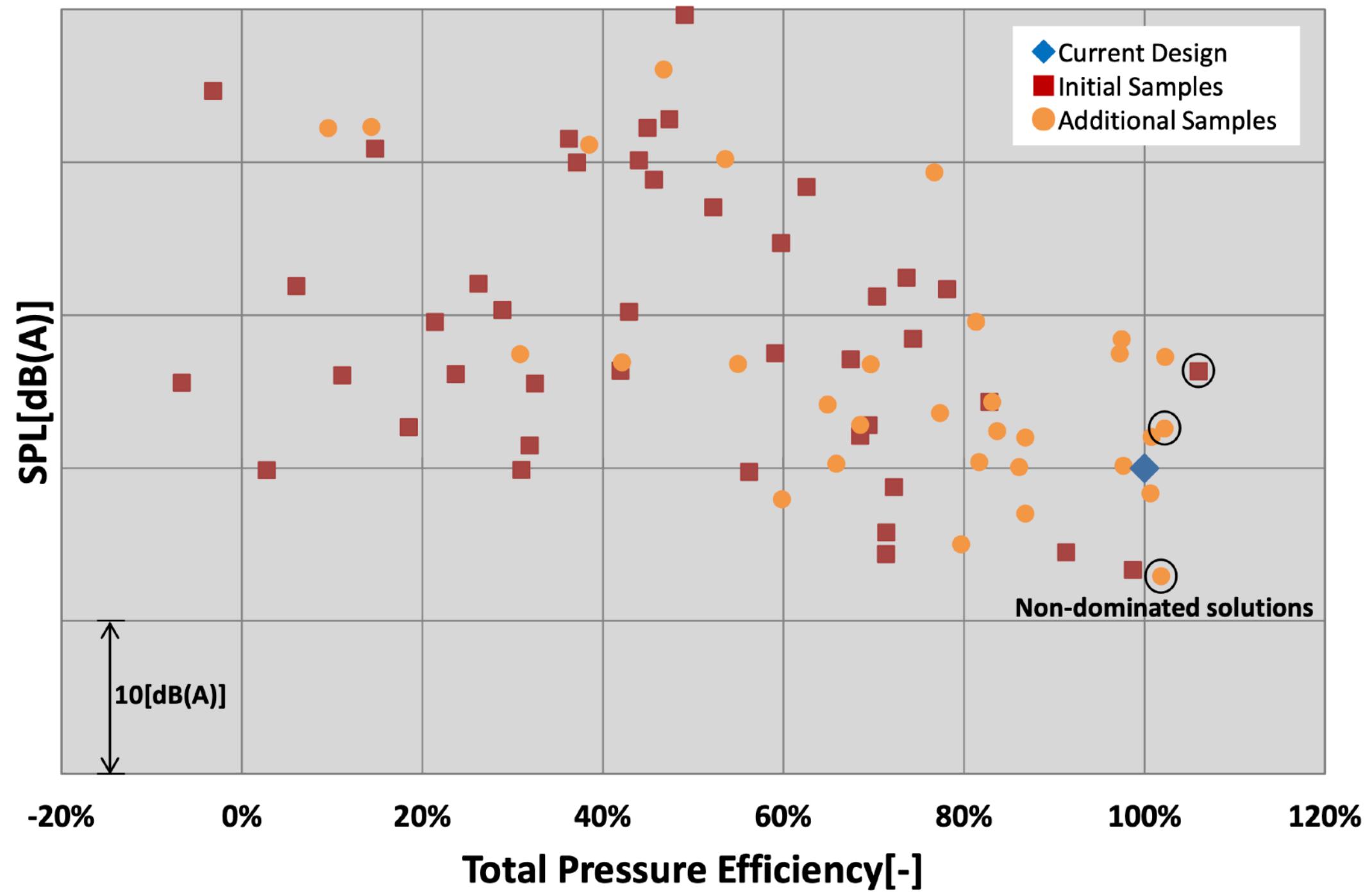
# Practical Example 2 : Car air conditioner design

Masaru Kamada, Koji Shimoyama, Fumito Sato, Junya Washiashi, and Yasufumi Konishi. Multi-objective design optimization of a high efficiency and low noise blower unit of a car air-conditioner. Proceedings of the Institution of Mechanical Engineers, Part D: Journal of Automobile Engineering, 233 (13):3493–3503, 2019.



- Determine the optimal design that maximizes airflow performance while minimizing noise.
- Perform fluid simulations → Evaluate performance
- 15 explanatory variables, such as the length and angle of the blades.

# Practical Example 2 : car air conditioner design



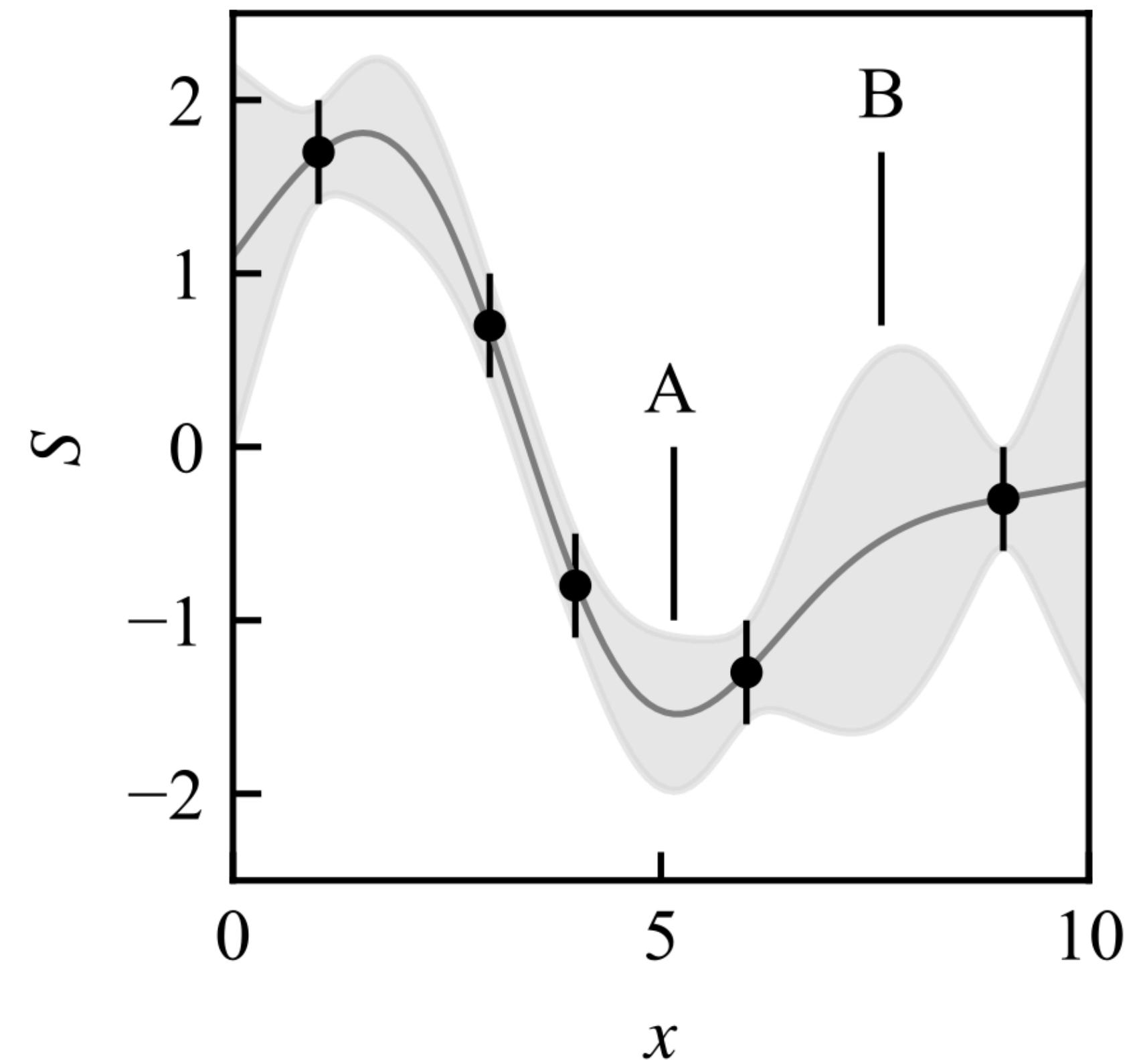
- Bayesian Optimization
  1. Simulate 45 different sets of design (= 15 explanatory variables) and measure the airflow performance and noise.
  2. Gaussian process regression across the 15-dimensional space of the explanatory variables.
  3. Based on EI (Expected Improvement), simulate the next 4 sets of designs with high airflow performance and low noise.
  4. Repeat this process 8 times.

$$EI(\boldsymbol{x}) = \int_{-\infty}^{y_{\min}} (y_{\min} - y(\boldsymbol{x})) \mathcal{N}(y_{\min} - \mu(\boldsymbol{x}), \sigma^2(\boldsymbol{x})) dy$$

# Hands-on exercise #3

Gaussian process regression

Lecture\_Day1\_Uemura/03\_GaussianProcess.ipynb



# **4. Classification models**

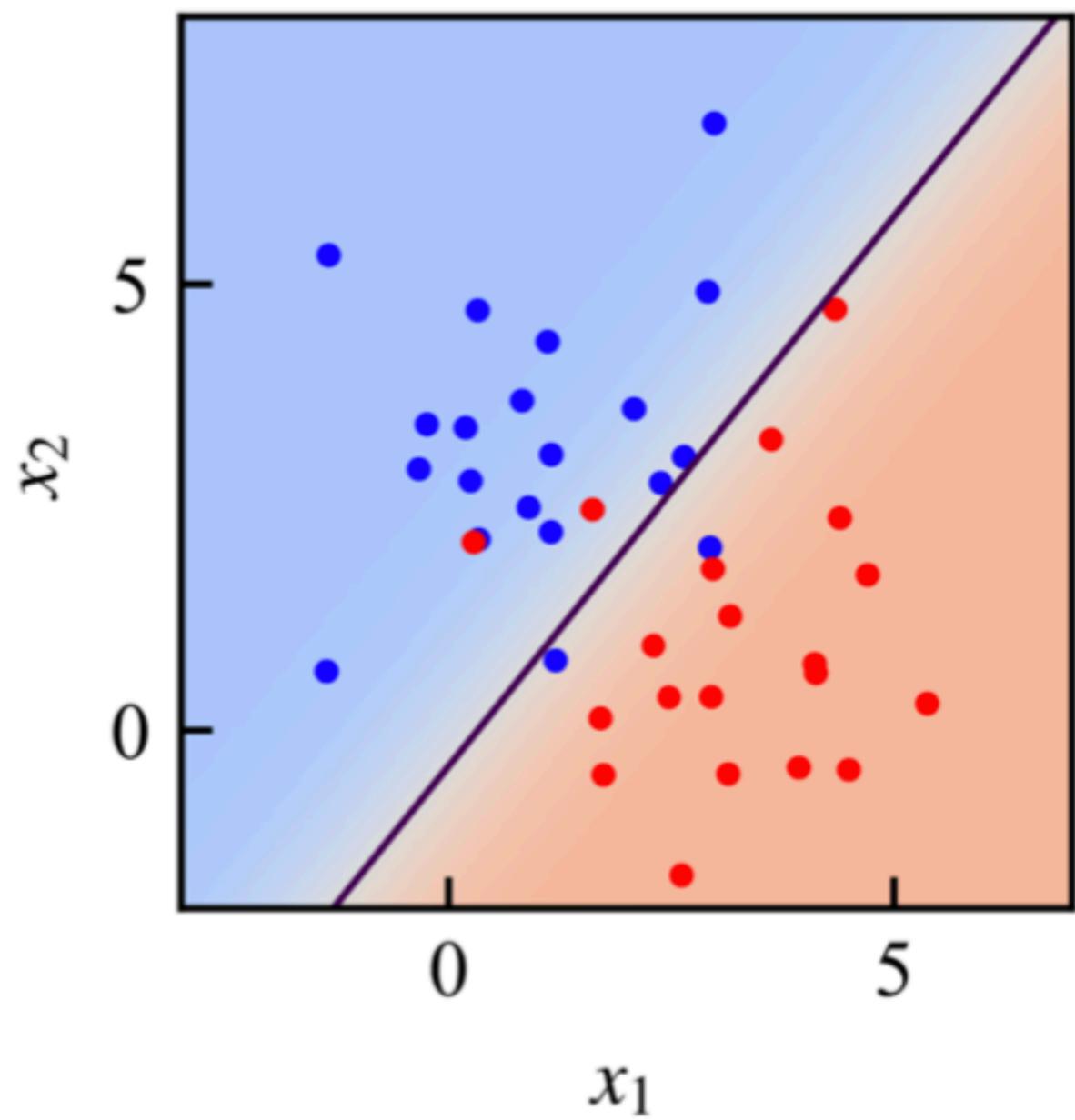
**4-1 Binary classification, evaluation metrics, confusion matrix**

**4-2 Models: Logistic regression & Support Vector Machine (SVM)**

## **4-1 Binary classification, evaluation metrics, confusion matrix**

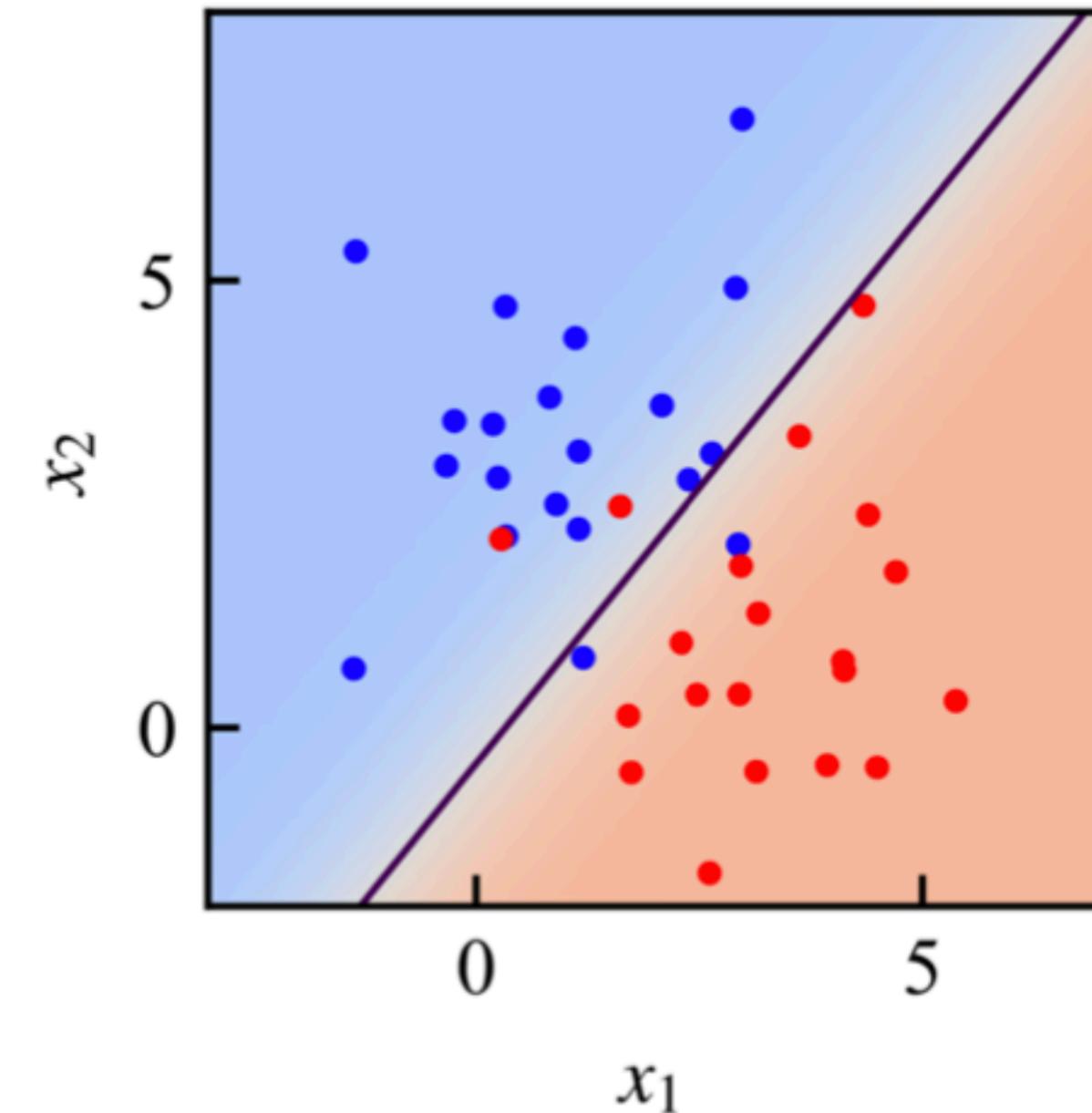
# Binary classification problems

- Binary classification: “Positive” or “Negative,” “Real” or “Fake,” or “Blue” or “Red”
  - While multi-class classification is an important topic, we won’t be covering it in this lecture due to time constraints.
- The target variable  $y$  has two possible classes,  $C_A$  and  $C_B$ ,
  - In SVM (as discussed later), denoted as,  $y = +1$  for  $C_A$  and  $y = -1$  for  $C_B$
- Determined by the explanatory variables  $\mathbf{x}$  (, also known as **the feature vector**).
- Develop a classification model:  $p(C_A | \mathbf{x})$ ,  $y=f(\mathbf{x})$ 
  - Maximizing or minimizing an evaluation metric.
  - While in regression problems we often use an evaluation function like  $\|\mathbf{y} - f(\mathbf{x})\|_2^2$  , in classification, we use different metrics.



# Error matrix

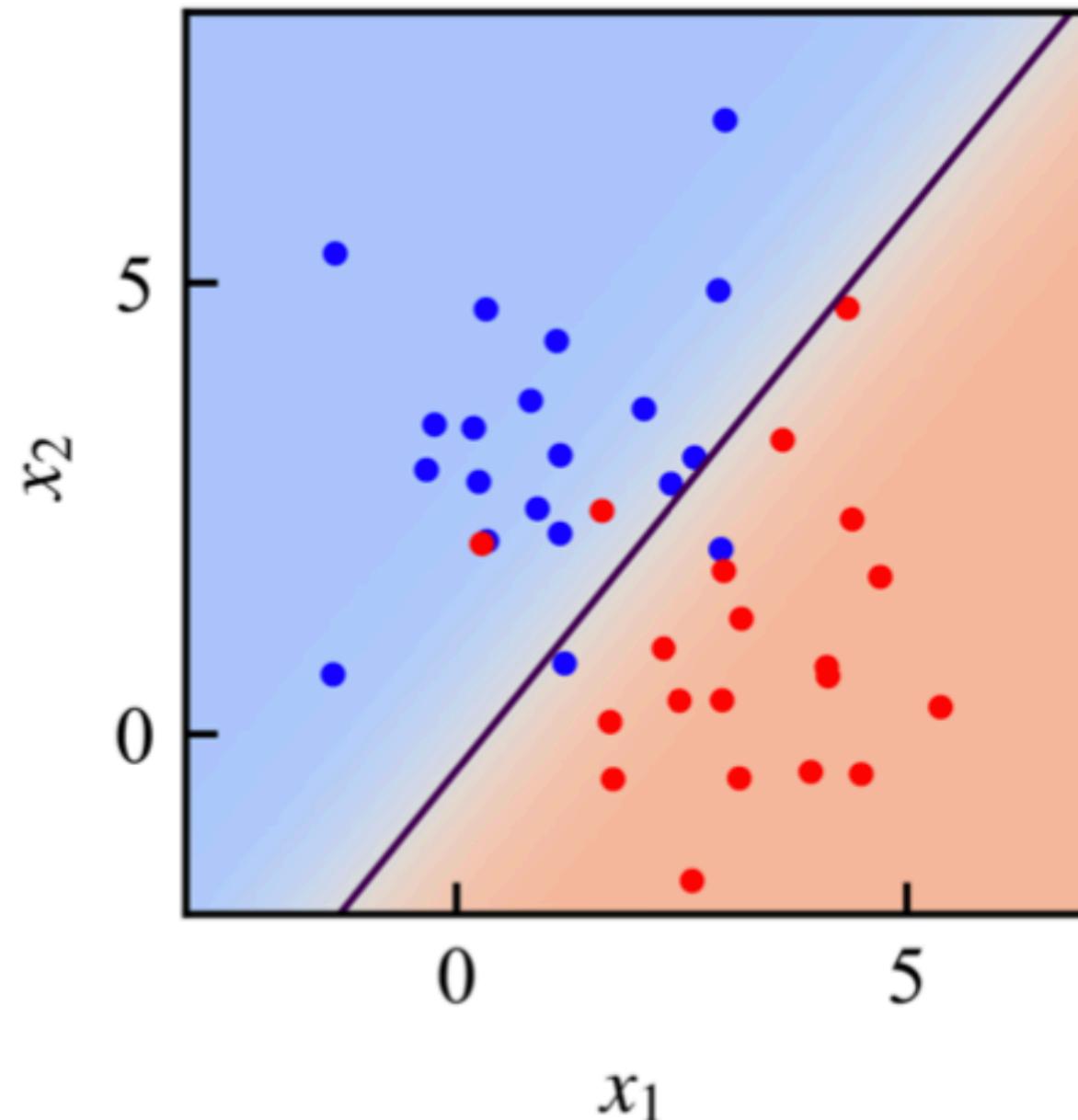
- Figure: An Example of a binary Classification
  - Two features
- Table: Error matrix (confusion matrix)
  - Of the 20 samples that are truly Class A:
    - 18 are correctly classified as Class A by the model (True Positive, TP)
    - 2 are incorrectly classified as Class B (False Negative, FN)
  - Of the 20 samples that are truly Class B:
    - 2 are incorrectly classified as Class A by the model (False Positive, FP)
    - 18 are correctly classified as Class B (True Negative, TN)
  - Hereafter, when we refer to TP, we are referring to the number of TP samples (ex. TP=18).



|           |       | True classes<br>(label) |            |
|-----------|-------|-------------------------|------------|
|           |       | $C_A$                   | $C_B$      |
| Estimates | $C_A$ | 18<br>(TP)              | 2<br>(FP)  |
|           | $C_B$ | 2<br>(FN)               | 18<br>(TN) |

# Evaluation metrics based on the error matrix

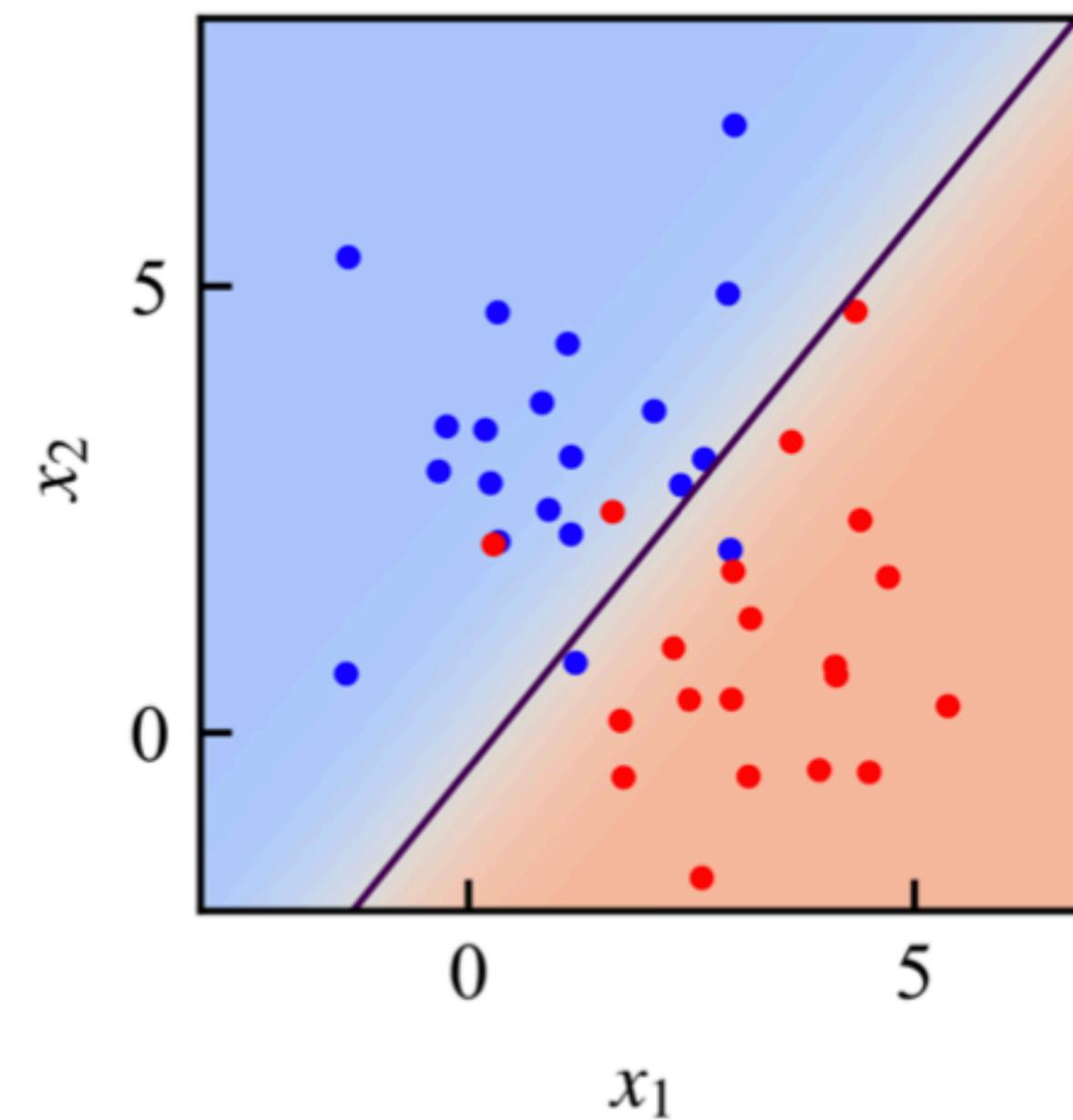
- Accuracy
  - The proportion of samples correctly classified out of all samples
  - $(TP+TN)/(Total\ number\ of\ samples) : (18+18)/40 = 0.9$
- Recall (True Positive Rate: TPR)
  - The proportion of samples correctly classified as Class A out of class A samples.
  - $TP/(TP+FN) : 18/(18+2)=0.9$
- False Positive Rate (FPR)
  - The proportion of samples incorrectly classified as Class A when they actually belong to Class B
  - $FP/(FP+TN) : 2/(18+2)=0.1$



|           |       | True classes<br>(label) |            |
|-----------|-------|-------------------------|------------|
|           |       | $C_A$                   | $C_B$      |
| Estimates | $C_A$ | 18<br>(TP)              | 2<br>(FP)  |
|           | $C_B$ | 2<br>(FN)               | 18<br>(TN) |

# Problems with Unequal Sample Sizes across Classes

- Ideally, we would have 20 samples each for class A and class B.
- What if we have 10,000 samples for class A and only 10 for class B?
- A model that classifies all samples as class A:
  - Accuracy:  $10000/10010 \sim 0.999$  & TPR = 1 → Looks like a good model
  - FPR =  $10/10 = 1$  → a significant problem
- When there is a significant imbalance in class sizes, it's crucial to consider both TPR and FPR.

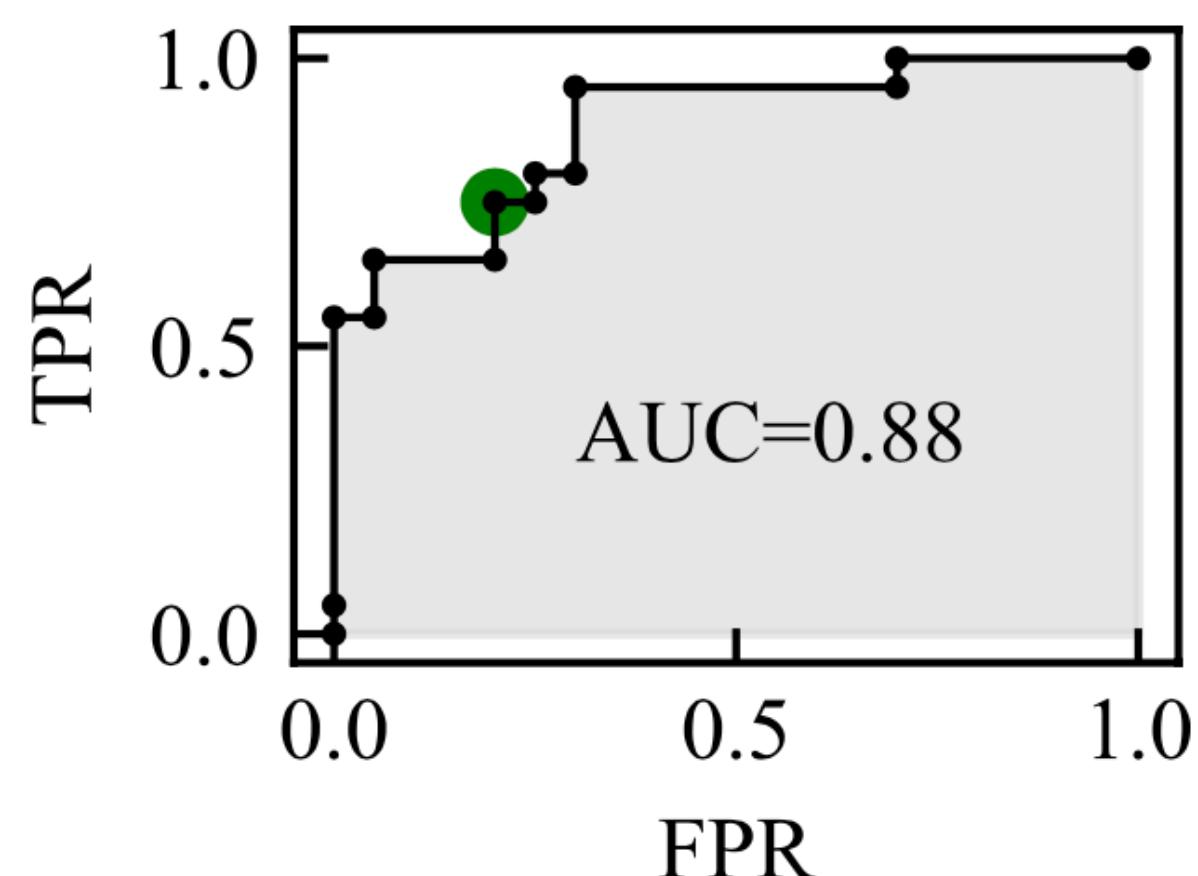
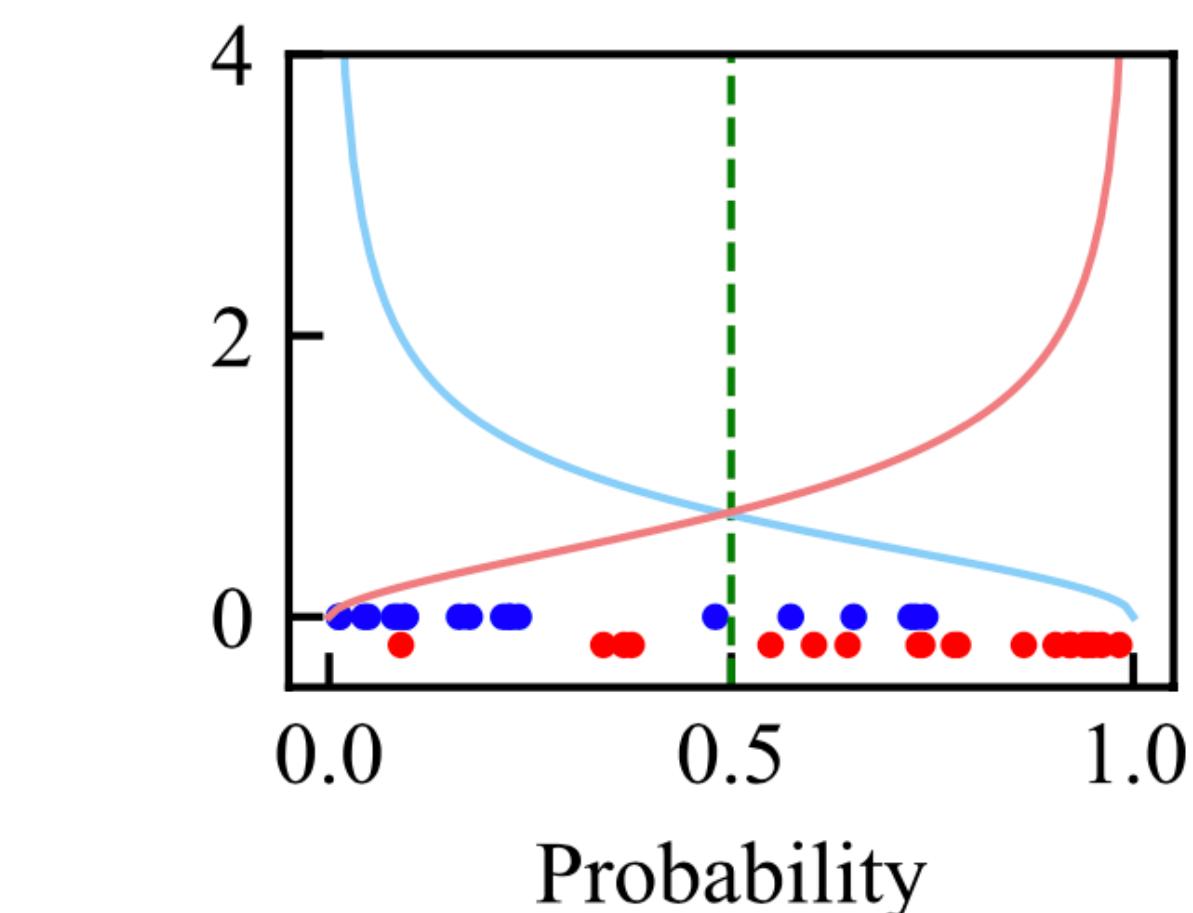
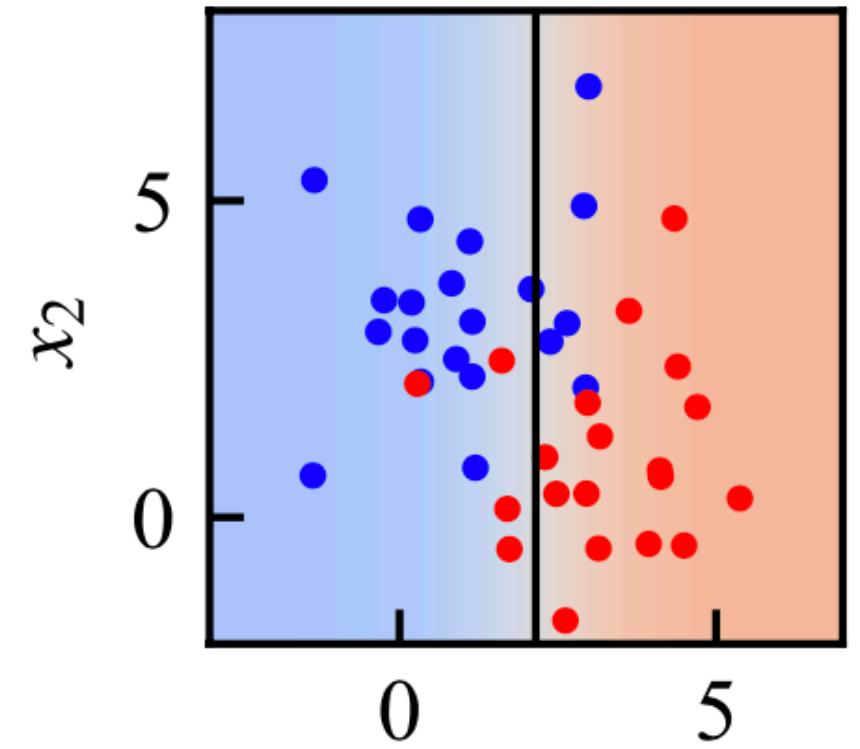
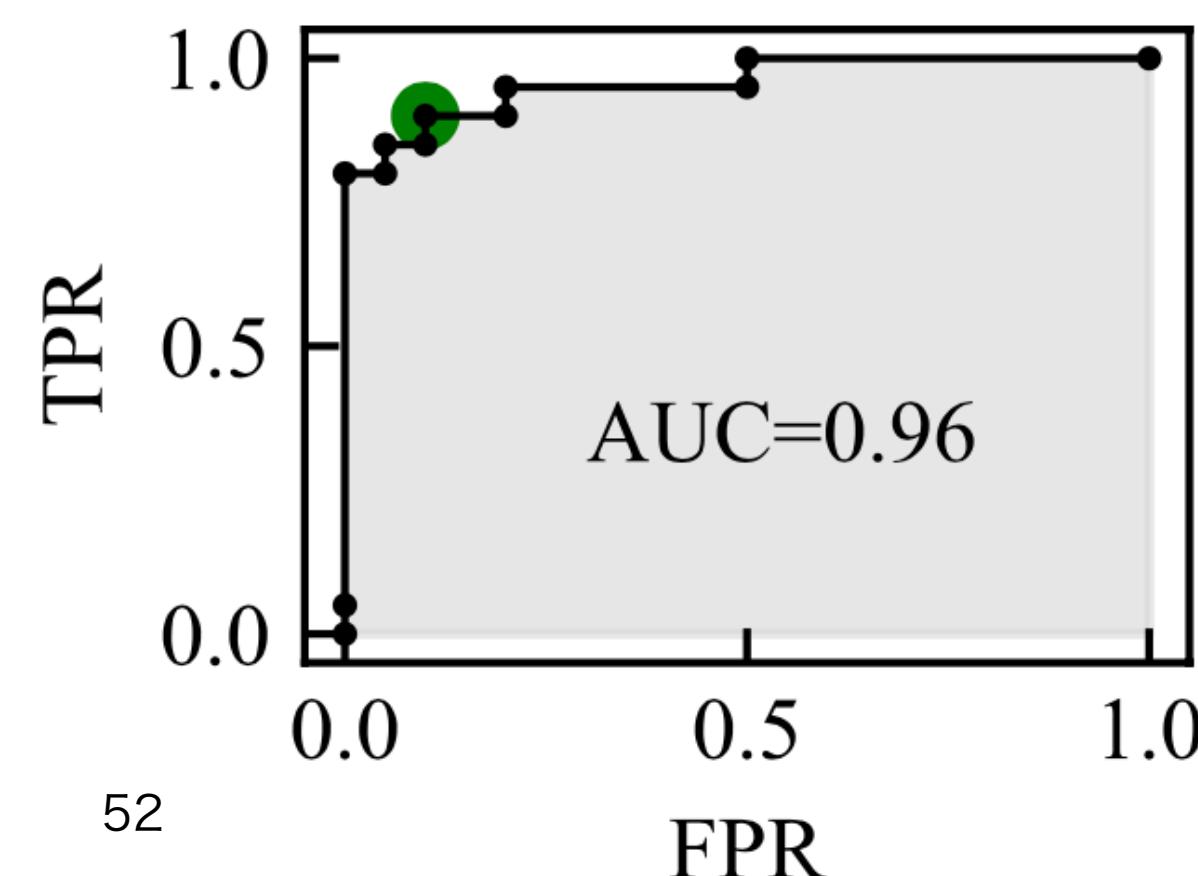
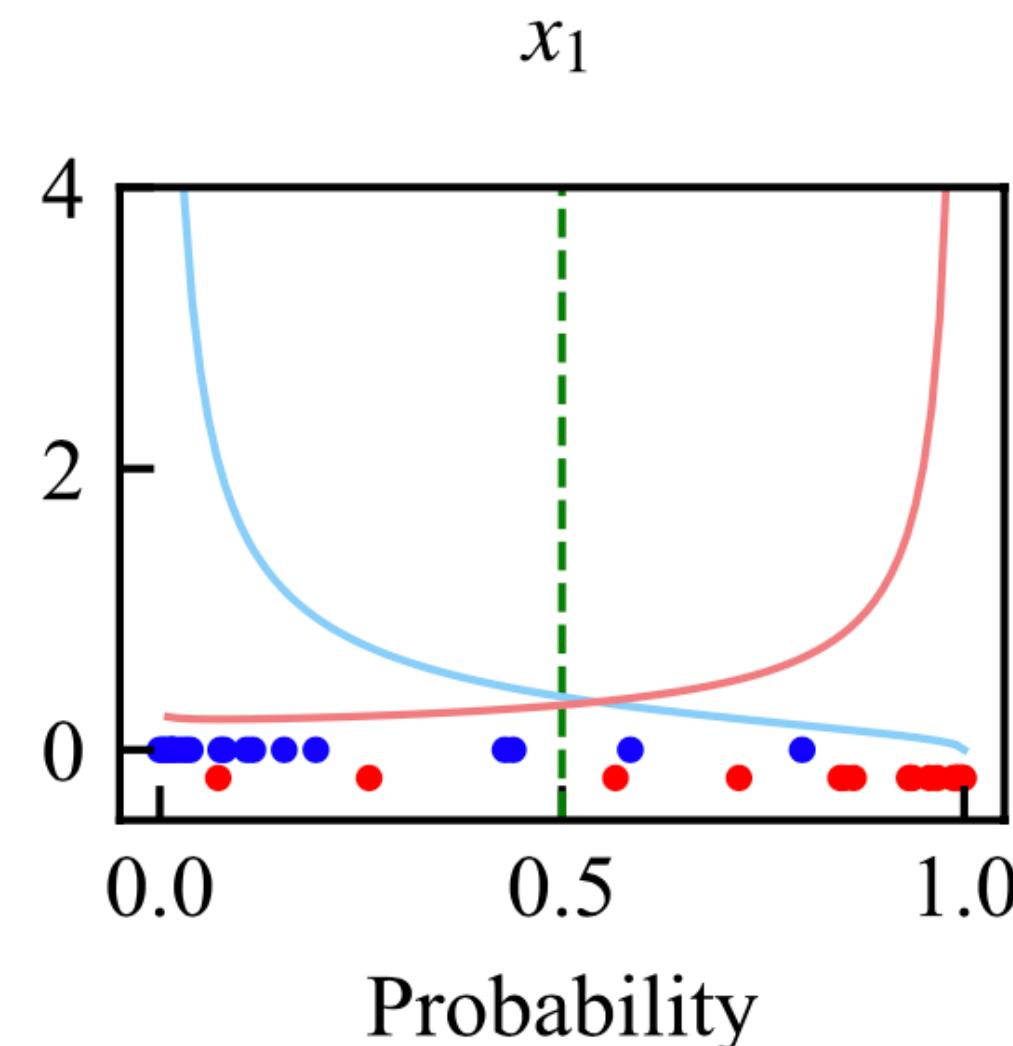
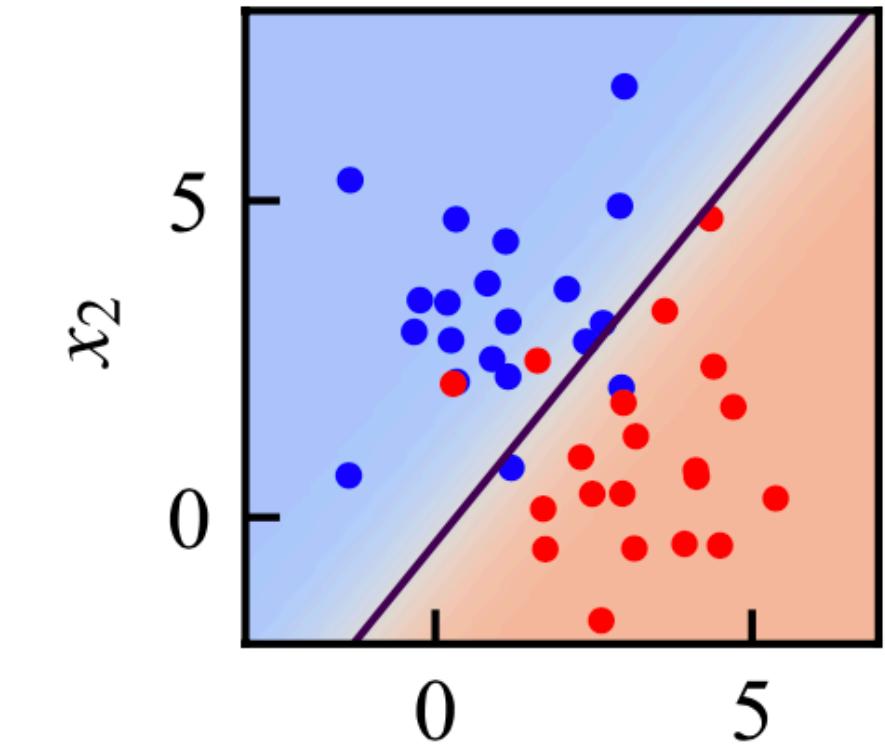


|           |       | True classes<br>(label) |            |
|-----------|-------|-------------------------|------------|
|           |       | $C_A$                   | $C_B$      |
| Estimates | $C_A$ | 18<br>(TP)              | 2<br>(FP)  |
|           | $C_B$ | 2<br>(FN)               | 18<br>(TN) |

# Receiver Operating Characteristic (ROC) Curve

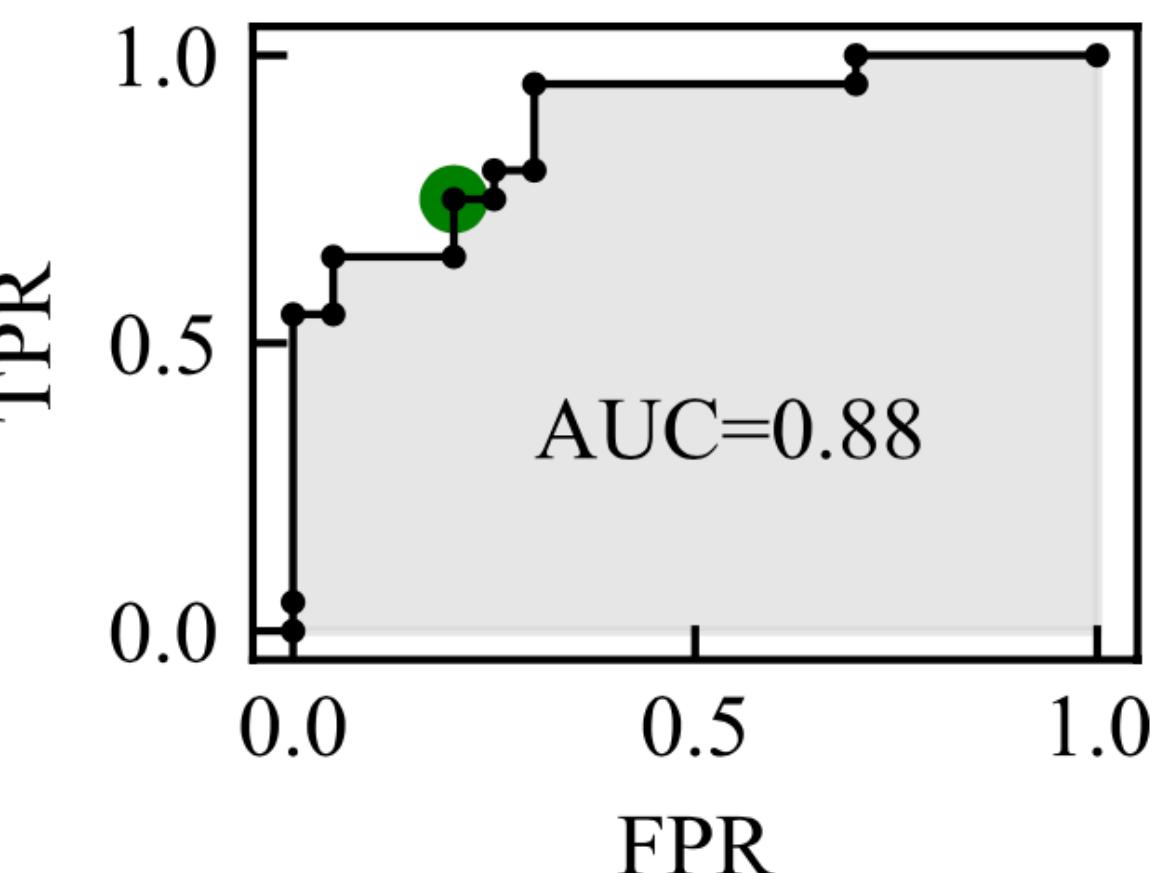
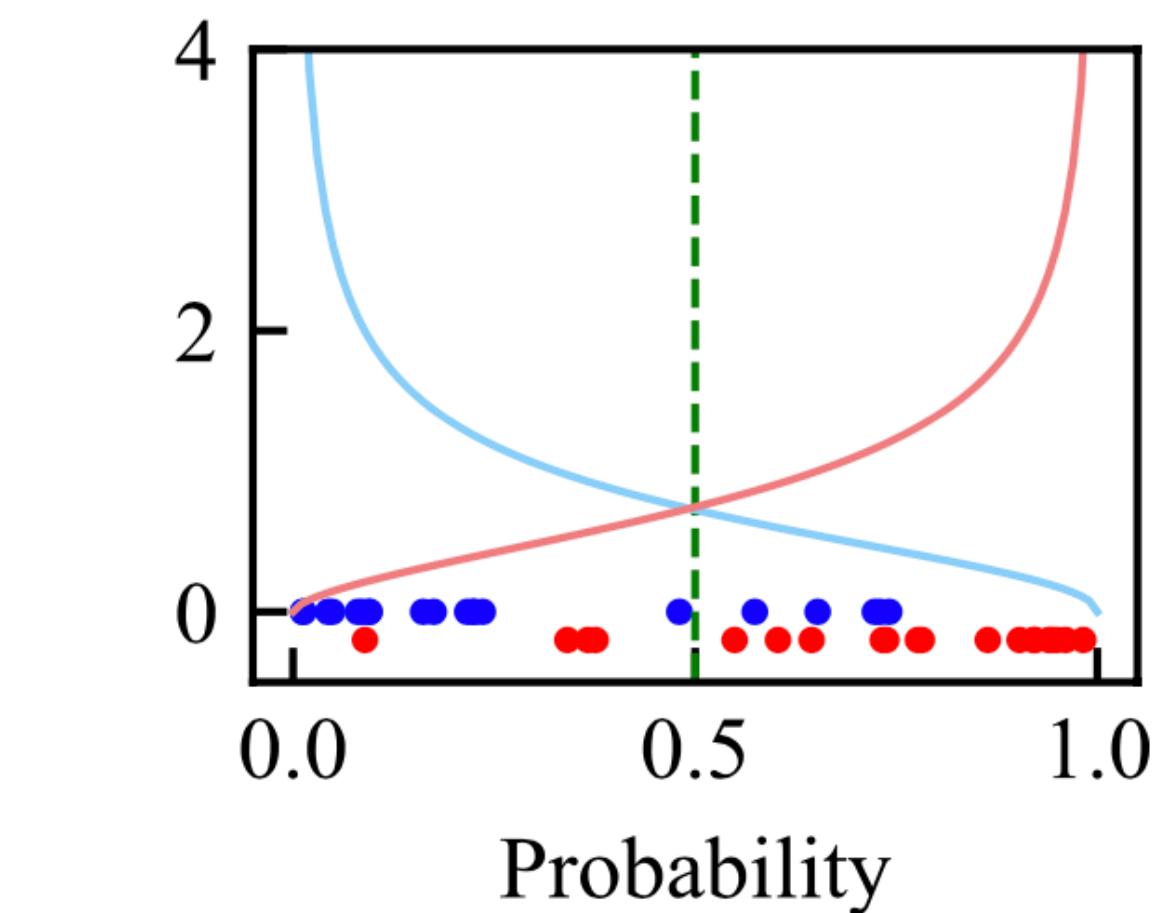
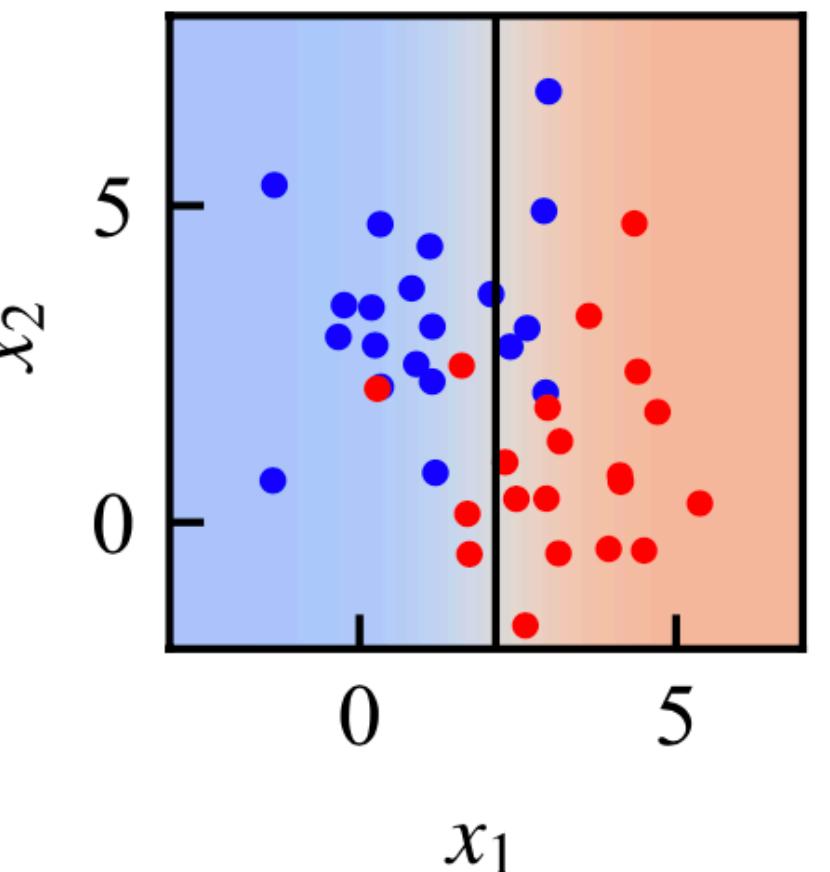
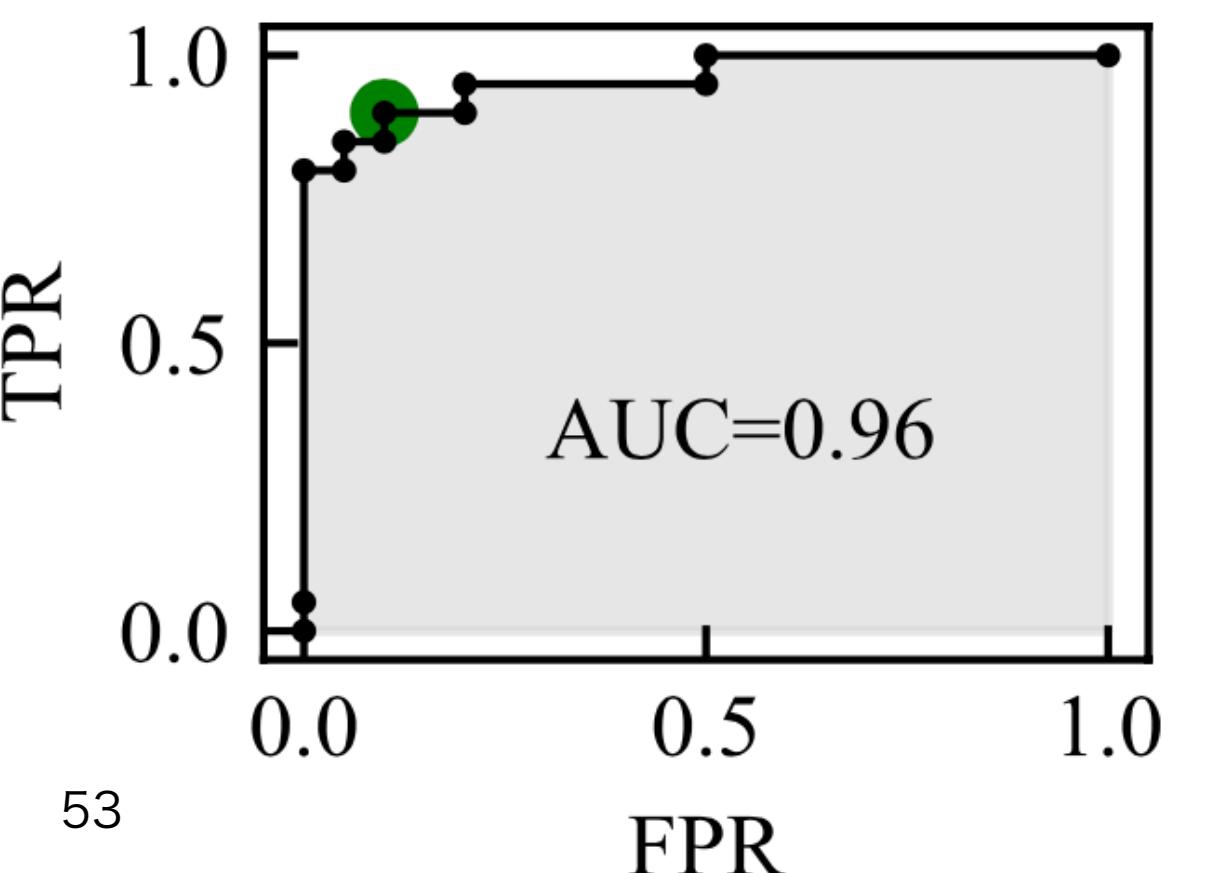
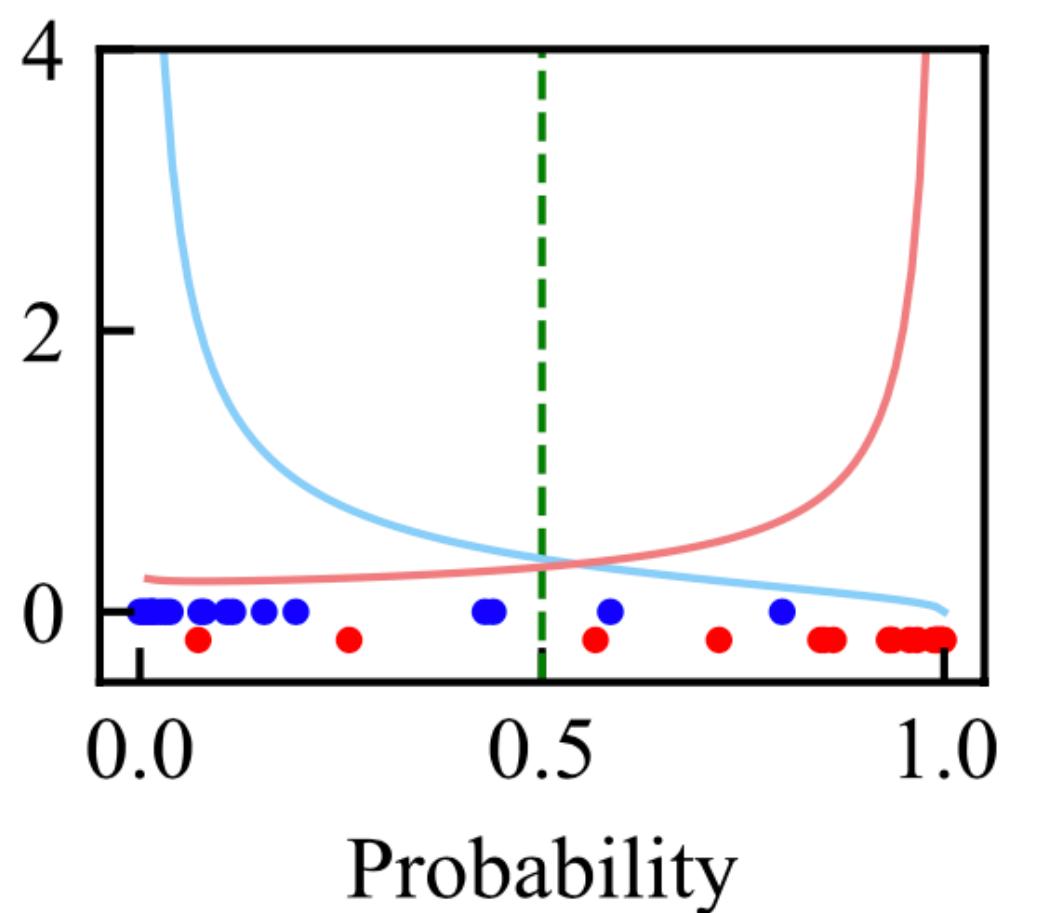
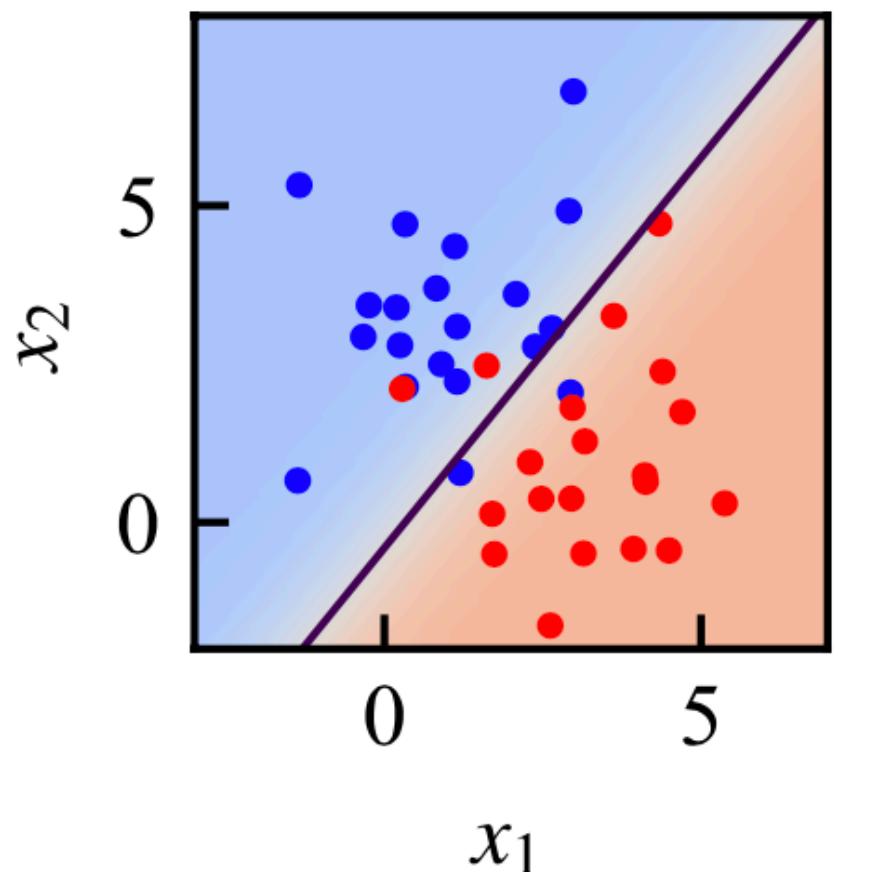
- Example on the right:

- With  $p=0.5$  as the decision boundary:  
 $\text{TPR}=18/20=0.9$ ,  $\text{FPR}=2/20=0.1$
- On the bottom left plot: This point is represented by a green dot, with FPR on the x-axis and TPR on the y-axis.
- With  $p=0.2$  as the decision boundary:  
 $\text{TPR}=19/20=0.95$ ,  $\text{FPR}=4/20=0.2$
- On the bottom left plot: The green dot moves to a point further up and to the right.
- By varying the threshold probability used for decision-making, you can plot a curve on the TPR-FPR plane  $\rightarrow$  the ROC curve.



# Evaluating ROC Curves and AUC

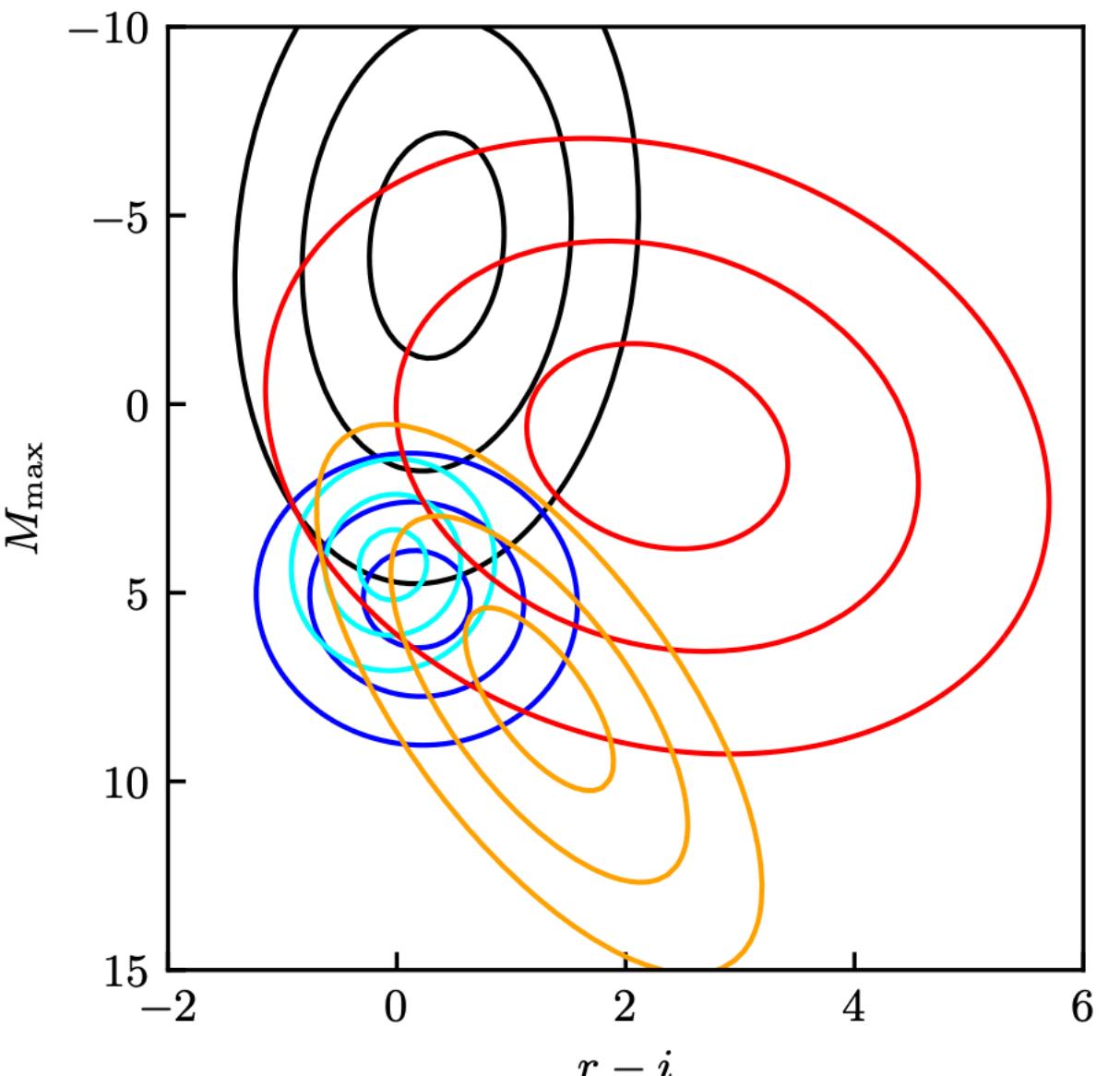
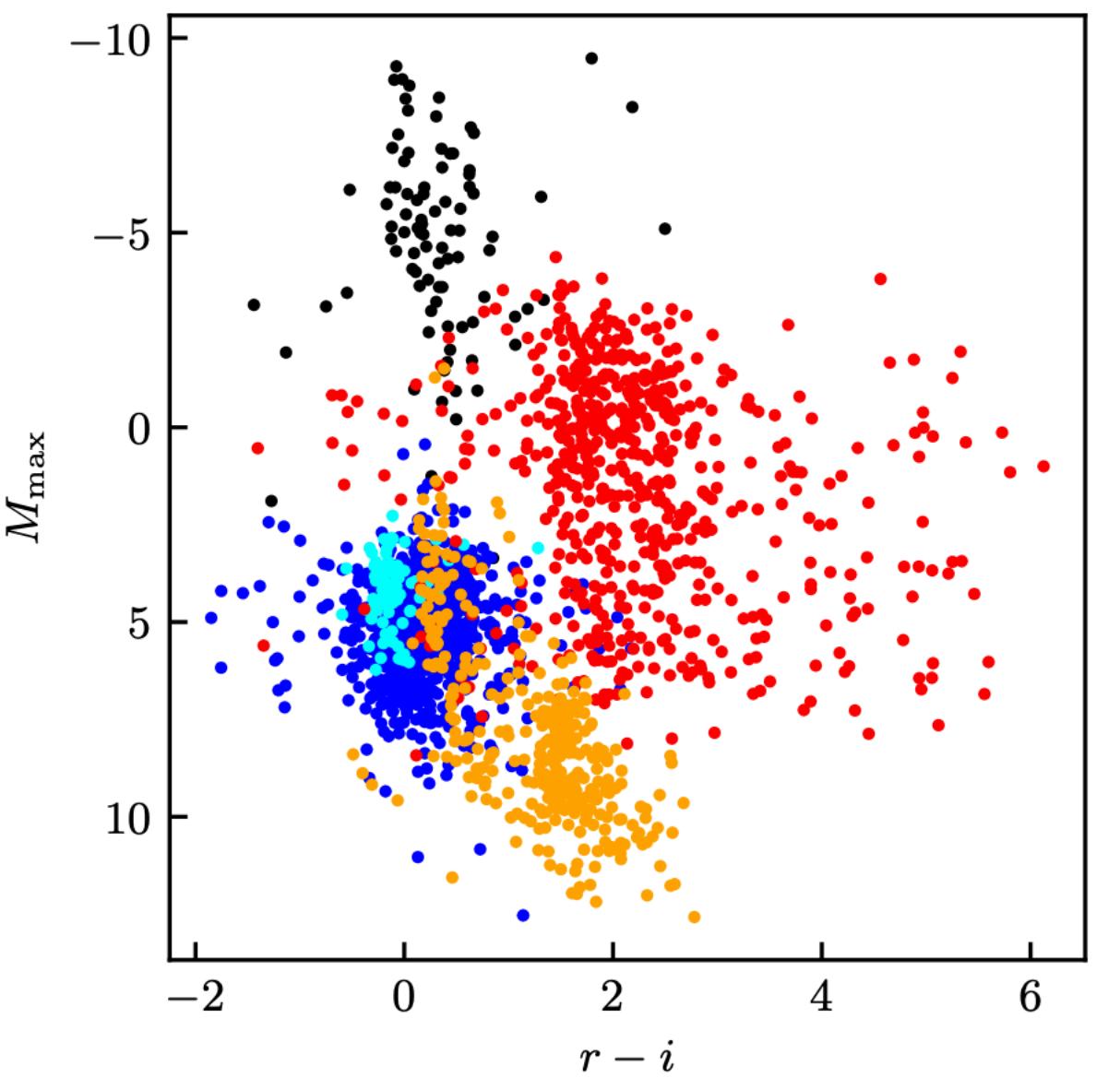
- An ideal classification model has  $\text{FPR}=0$  &  $\text{TPR}=1$
- The proximity of the ROC curve to the point  $(0,1)$  indicates the model's performance.
- Example in the right panels:
  - The model using only  $x_1$  results in an ROC curve that is farther from  $(0,1)$   $\rightarrow$  Indicates a lower-performance model.
- The Area Under the ROC Curve (AUC) is a commonly used performance metric:
  - Larger the AUC, Better the model.
  - Left panel:  $\text{AUC} = 0.96$ , Right panel:  $\text{AUC} = 0.88$   $\rightarrow$  The model in the left figure performs better.



# **4-2 Models: Logistic regression & Support Vector Machine (SVM)**

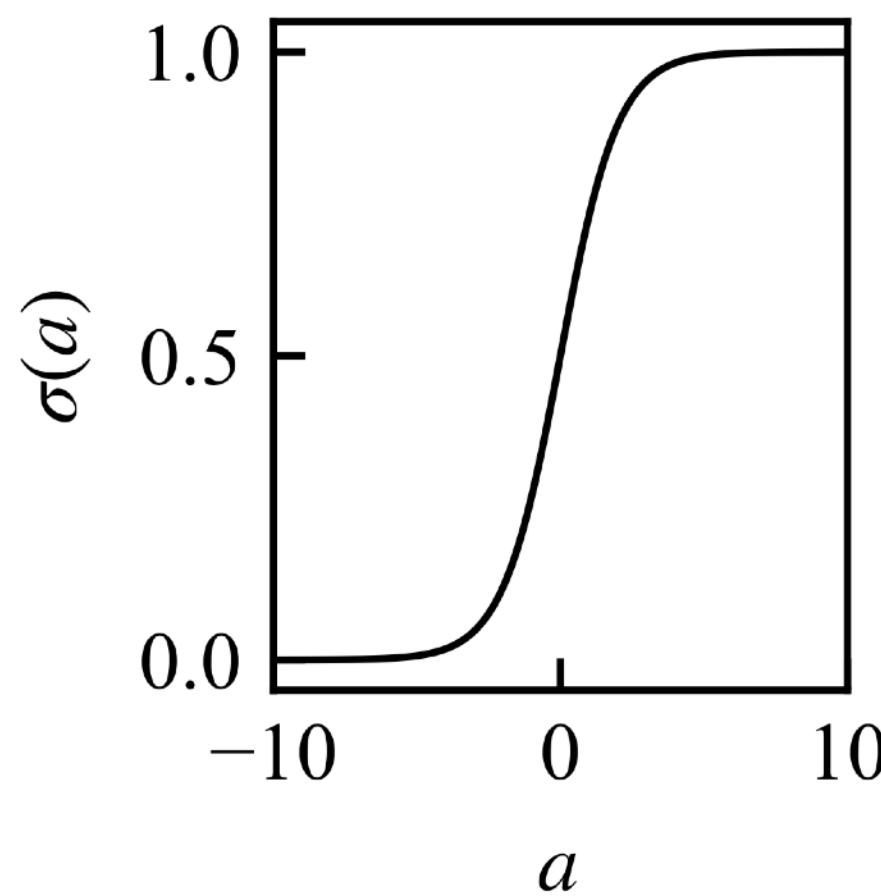
# Probabilistic Classification Models

- The probability that a sample with feature vector  $\mathbf{x}$  belongs to  $y = C_A$  is given by  
$$p(y = C_A | \mathbf{x}) = p(C_A | \mathbf{x}) .$$
- Using Bayes' theorem:  
$$\begin{aligned} p(C_A | \mathbf{x}) &= \frac{p(\mathbf{x}|C_A)p(C_A)}{\sum p(\mathbf{x}|y)p(y)} \\ &= \frac{p(\mathbf{x}|C_A)p(C_A)}{p(\mathbf{x}|C_A)p(C_A) + p(\mathbf{x}|C_B)p(C_B)} \end{aligned}$$
- Bayesian classifier: A model that computes the posterior probability by providing all elements on the right side of Bayes' theorem.
  - $p(x|C_A)$  : The distribution of  $\mathbf{x}$  for samples belonging to class  $C_A$  (e.g., approximated by a Gaussian distribution).
  - $p(C_A)$  : The prior probability that a sample belongs to class  $C_A$  (e.g., the proportion of  $C_A$  samples in the training data).
  - Requires adjusting many parameters (e.g., the covariance matrix in the case of a multivariate normal distribution).



# Logistic regression

- $\sigma(a)$  : Logistic sigmoid function
- $a$  is a function of the feature vector  $\mathbf{x}$  (linear combination)
- Estimate  $\mathbf{w}$  (weights) from the data using maximum likelihood estimation = **Logistic Regression**.
- Target variable: Expressed as  $y = 1$  for class A and  $y = 0$  for class B.
- Negative log-likelihood :  $E(\mathbf{w})$ 
  - Takes the form of **cross-entropy** between  $\mathbf{y}$  and  $\sigma$ .



$$\begin{aligned} p(C_A|\mathbf{x}) &= \frac{p(\mathbf{x}|C_A)p(C_A)}{p(\mathbf{x}|C_A)p(C_A) + p(\mathbf{x}|C_B)p(C_B)} \\ &= \frac{1}{1 + \exp(-a)} = \sigma(a) \end{aligned}$$

$$a = \log \frac{p(\mathbf{x}|C_A)p(C_A)}{p(\mathbf{x}|C_B)p(C_B)}$$

$$a = \mathbf{w}_1^T \mathbf{x} + w_0 = \mathbf{w}^T \mathbf{x}$$

$$p(y_i|\mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x}_i)$$

$$p(\mathbf{y}|\mathbf{w}) = \prod_i \sigma(\mathbf{w}^T \mathbf{x}_i)^{y_i} \{1 - \sigma(\mathbf{w}^T \mathbf{x}_i)\}^{1-y_i}$$

$$E(\mathbf{w}) = -\log p(\mathbf{y}|\mathbf{w}) = -\sum_i \{y_i \log \sigma(\mathbf{w}^T \mathbf{x}_i) + (1 - y_i) \log(1 - \sigma(\mathbf{w}^T \mathbf{x}_i))\}$$

# Deterministic Classification Models

- Probabilistic models = based on Bayes' theorem.
- Deterministic models construct a discriminant function  $f(x; \theta)$  and a decision boundary  $f(x; \theta) = c$ .
- There are various models, but here we introduce the representative **Support Vector Machine (SVM)**.

# Hard-Margin SVM (Perfectly Separable Problems)

- Linear discriminant function for features :

$$f(x) = w^T x + b$$

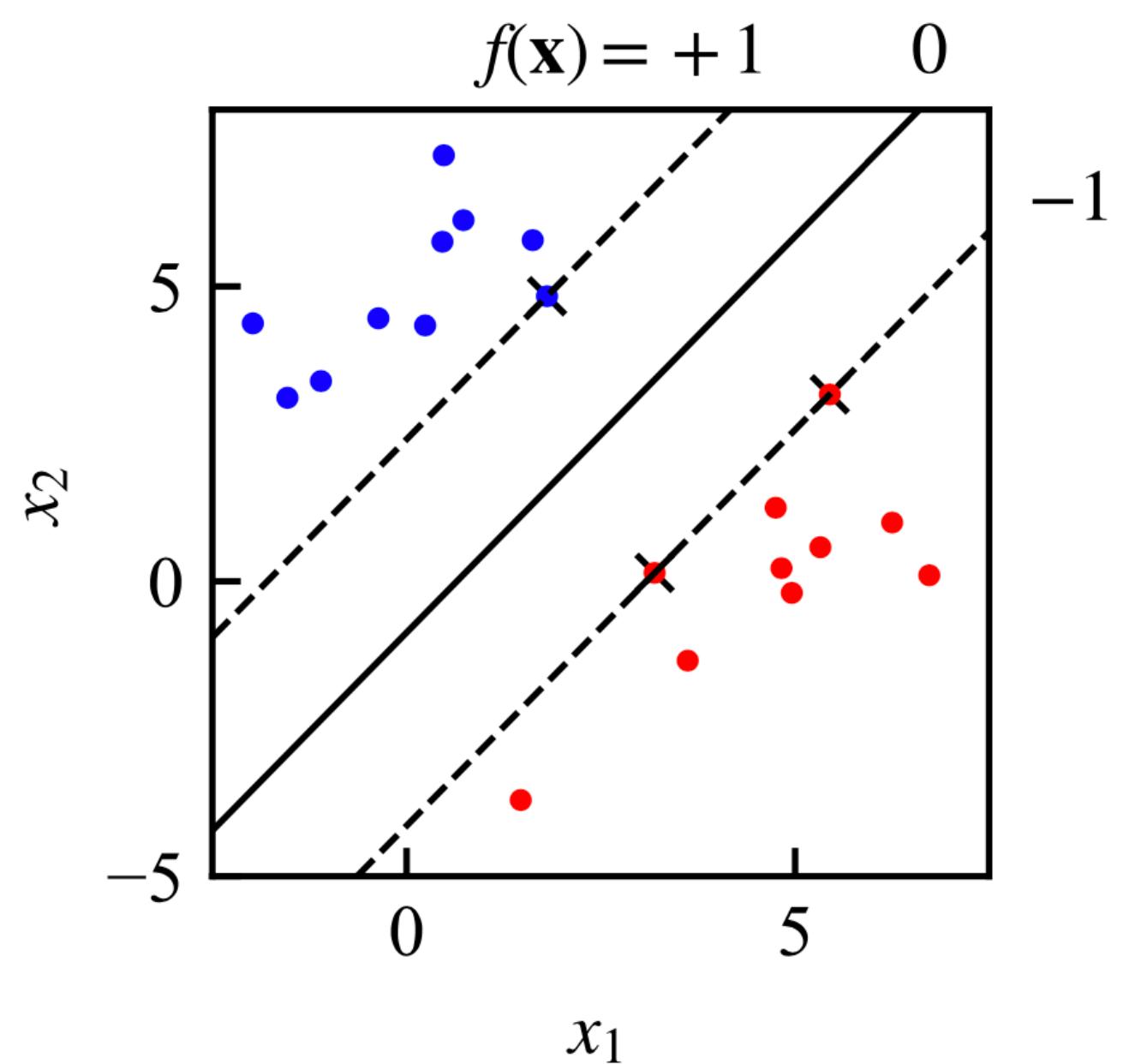
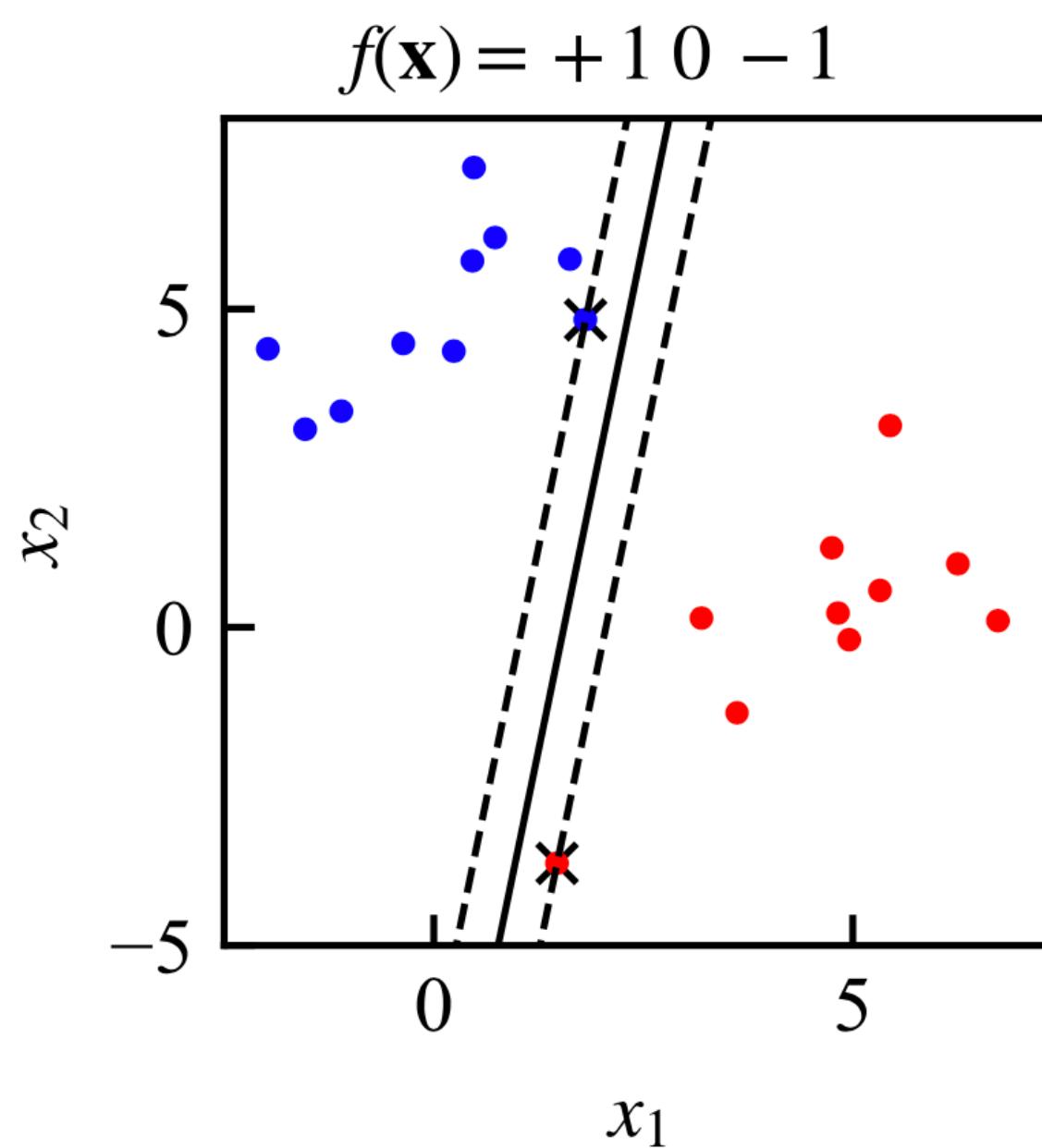
- Later, we'll extend this to nonlinear discriminant functions.

- Decision Boundary:  $f(x) = 0$

- Target Variable:  $y=1$  if class A,  $y=-1$  if class B

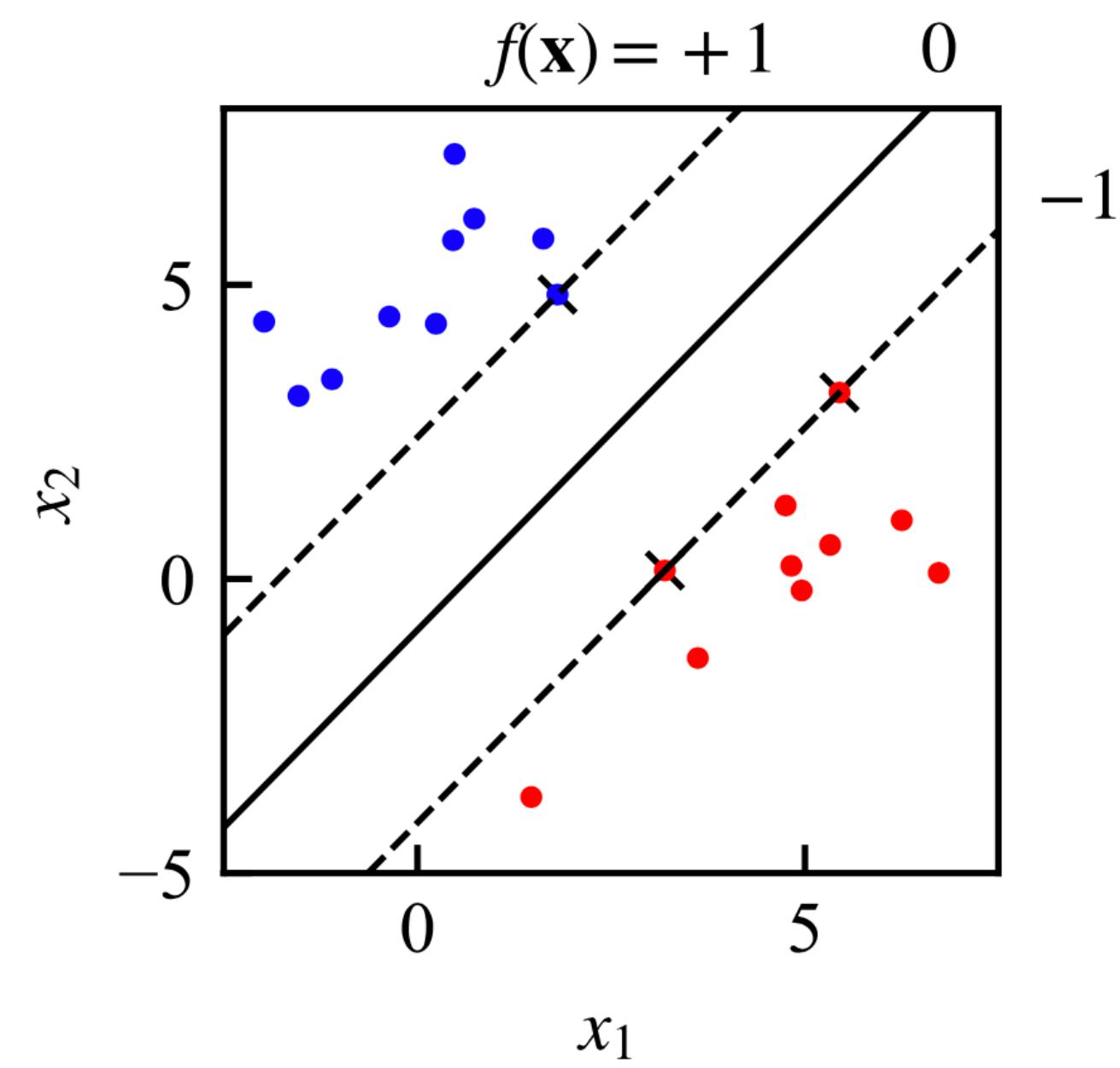
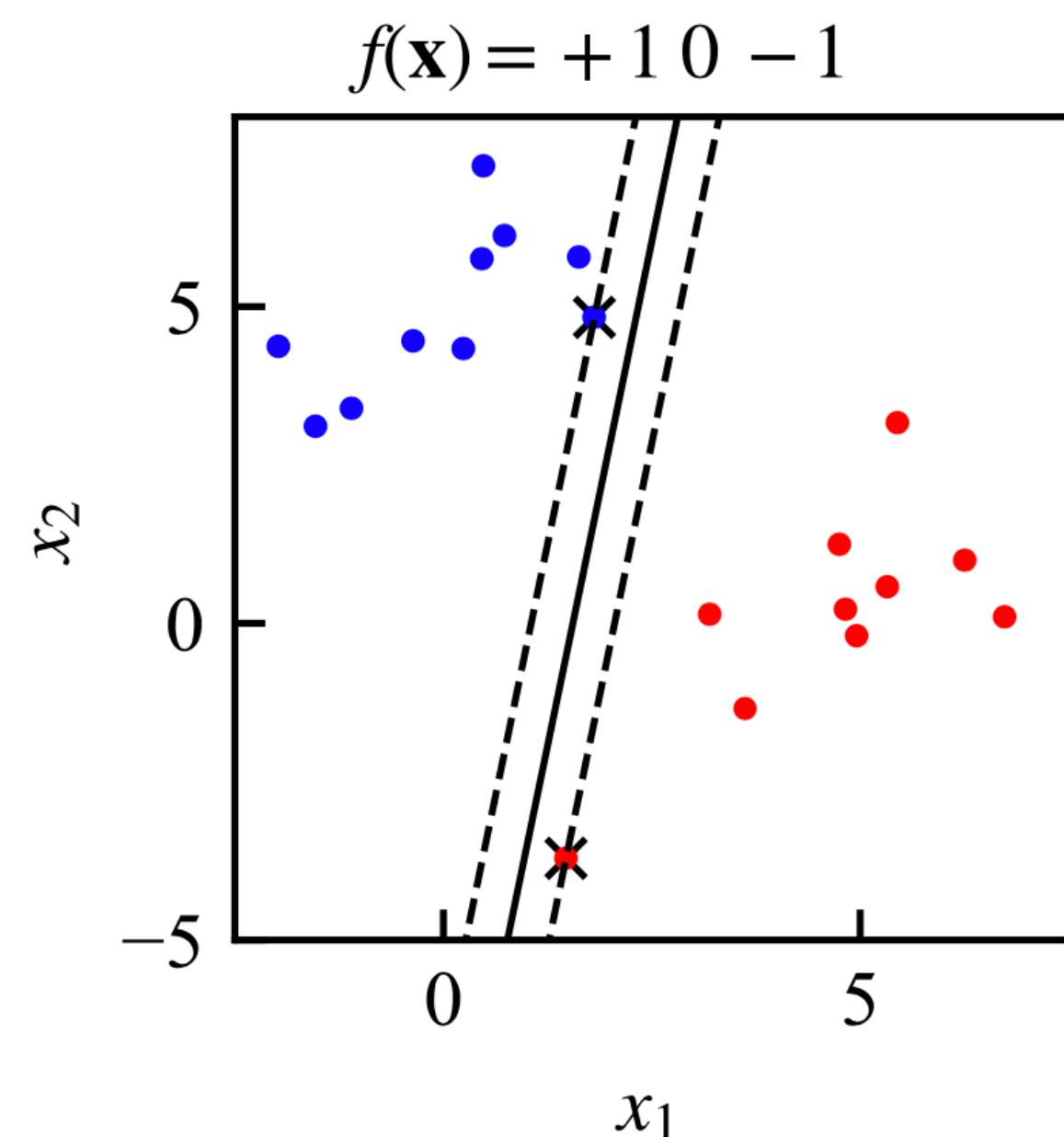
- The model classifies a sample as class A if  $f(x) > 0$  , class B if  $f(x) < 0$  .

- As illustrated in the right figure, there are numerous possible discriminant functions.



# SVM Chooses the Model with the Maximum Margin

- Margin: The distance between the decision boundary and the closest samples  $x_{SV}$ .
  - , called **the support vectors**.
- In general, the distance between  $x'$  and the hyperplane  $f(x) = w^T x = 0$  is  $|f(x')|/\|w\|_2$
- If we construct the decision function such that  $f(x_{SV}) = \pm 1$ , the margin becomes  $1/\|w\|_2$
- For Class A ( $y=1$ ),  $f(x) > 1$ , and for Class B ( $y=-1$ ),  $f(x) < -1$ . Therefore, the SVM problem reduces to the quadratic programming problem, as shown on the right.



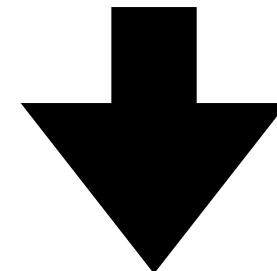
$$\begin{aligned} & \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2 \\ \text{subject to } & y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \end{aligned}$$

# Wanting a Non-Linear Boundary

- Apply a non-linear mapping  $\phi(x)$  ?
- Transform the quadratic programming problem into a single objective function using the Lagrange multipliers method.
  - Introduce new variables,  $\alpha$
- Minimize the objective function  $L$  by finding the conditions where the derivative of  $L$  is zero.
- Substitute the conditions for  $w$  and  $b$  back into the objective function.
- This transforms into a problem of maximizing  $L$  with respect to  $\alpha$ .
- The inner product of feature vectors  $x_i$  and  $x_j$  appears → Kernel trick.

$$\min_{w,b} \frac{1}{2} \|w\|_2$$

subject to  $y_i(w^T x_i + b) \geq 1$



$$L = \frac{1}{2} \|w\|_2 - \sum_i \alpha_i \{y_i(w^T x_i + b) - 1\}$$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w - \sum_i \alpha_i y_i x_i = 0$$

$$w = \sum_i \alpha_i y_i x_i$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow - \sum_i \alpha_i y_i = 0$$

$$\sum_i \alpha_i y_i = 0$$

$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i^T x_j$$

# The Kernel Trick

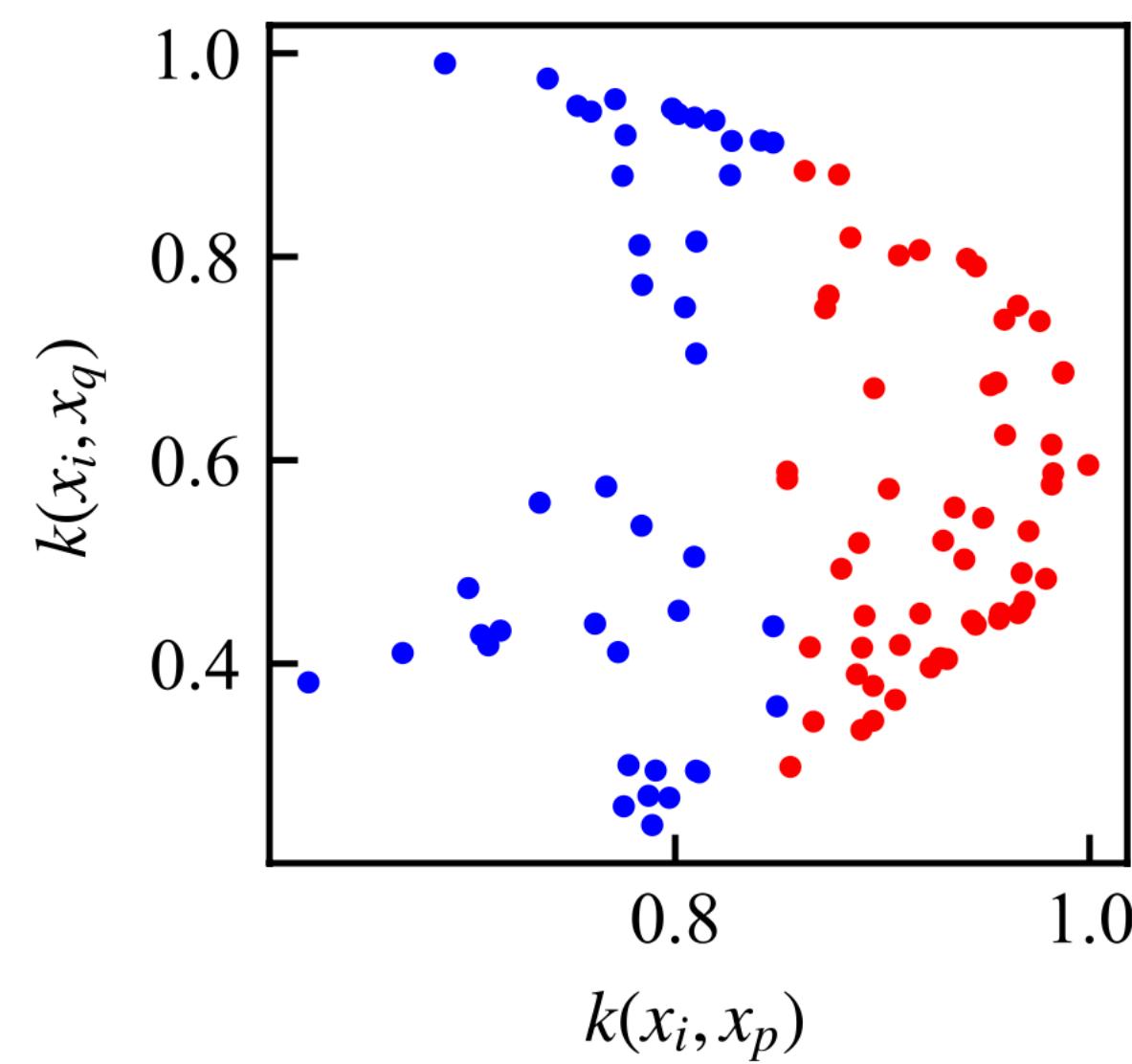
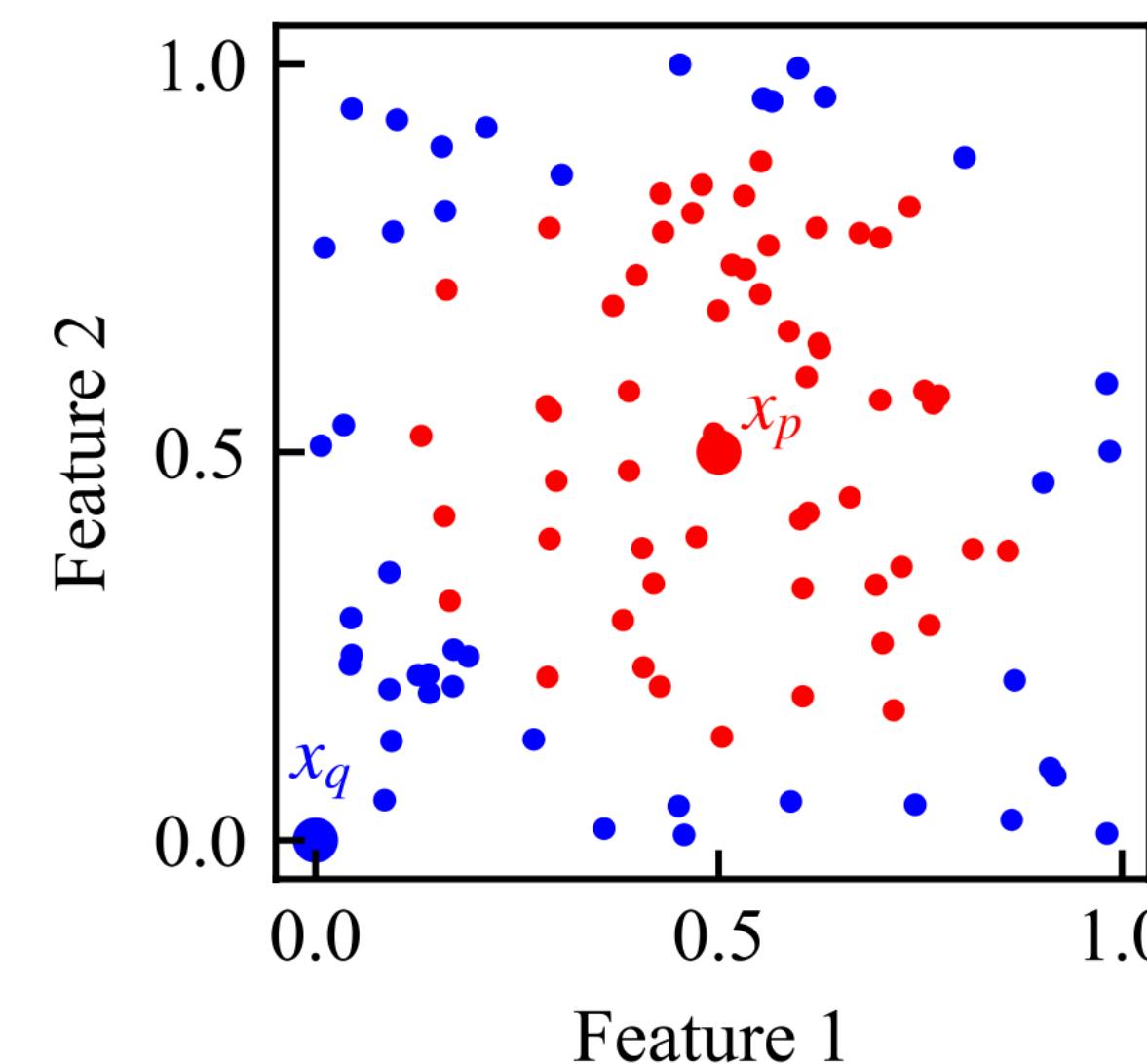
- **The Kernel Trick:** Replace the inner product of feature vectors with a kernel function which represent a measure of similarity between features.

- For example:

- the RBF (Radial Basis Function) kernel, or Gaussian kernel.
  - When  $x_i$  and  $x_j$  are close in feature space,  $k=1$ ; when they are far apart,  $k=0$ .
- Kernel functions typically have adjustable parameters (e.g.,  $\sigma$  for the RBF kernel).

$$L = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

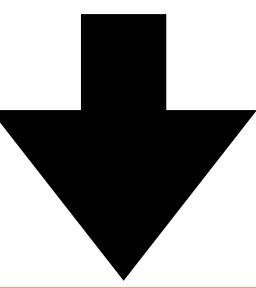


# Soft Margin SVM

- In real-world problems, it's rare to have perfectly separable data where hard margin SVM is applicable. Measurement errors and other influences often lead to inevitable misclassification → **Soft Margin SVM**.
- The right formula:
  - Hard Margin: For  $y = 1$ ,  $f(x) (=wx+b) > 1$  must hold.
  - Soft Margin: Relaxing the constraints → introducing new variables  $\xi$  ( $>0$ ), allowing  $f(x) < 1$  even for  $y = 1$
- If many  $\xi$  are non-zero, the performance of the classifier may degrade → add constraints on  $\xi$ .
- Adjust the regularization parameter  $C$  (using cross-validation, etc.) to build a model with high generalization performance.

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2$$

subject to  $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

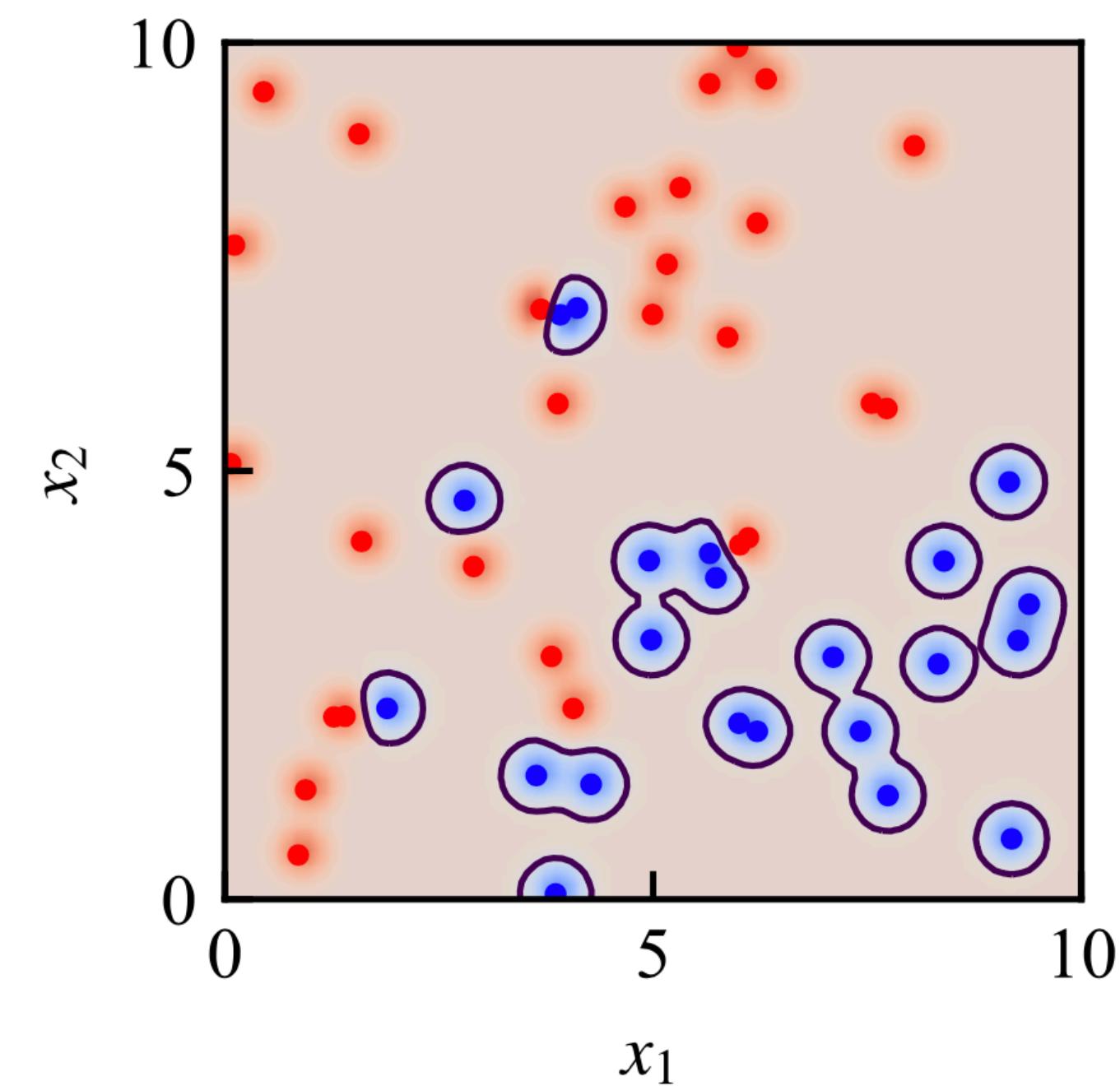
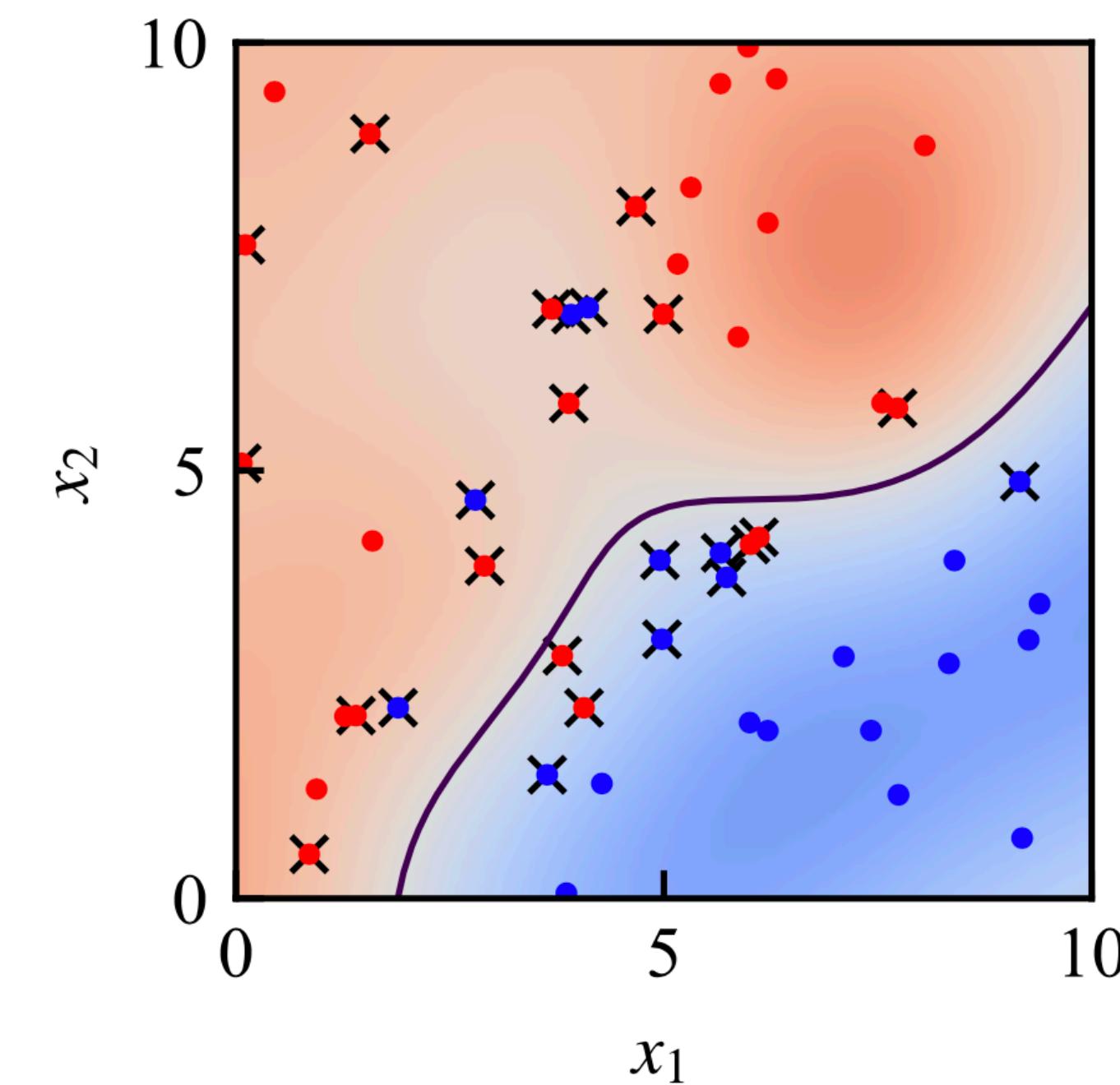


$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_i \xi_i$$

subject to  $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i$

# Classifying with Synthetic Data

- Use RBF kernel
- Determine **regularization** parameter C and **kernel** parameter  $\sigma$  by evaluating with **ROC-AUC** through **cross-validation**
- Left Figure: Best model
- Right Figure: **Overfitted** model
  - Kernel variance is small, and C is large.



# Hands-on exercise #4

Support Vector Machine

Lecture\_Day1\_Uemura/04\_SVM.ipynb

