

# Renormalized Mutual Information for Artificial Scientific Discovery

Leopoldo Sarra

*Max Planck Institute for the Science of Light  
Erlangen, Germany*

work with **Andrea Aiello** and **Florian Marquardt**

Learning to Discover - AI and Physics Conference  
Thursday, April 26th 2022



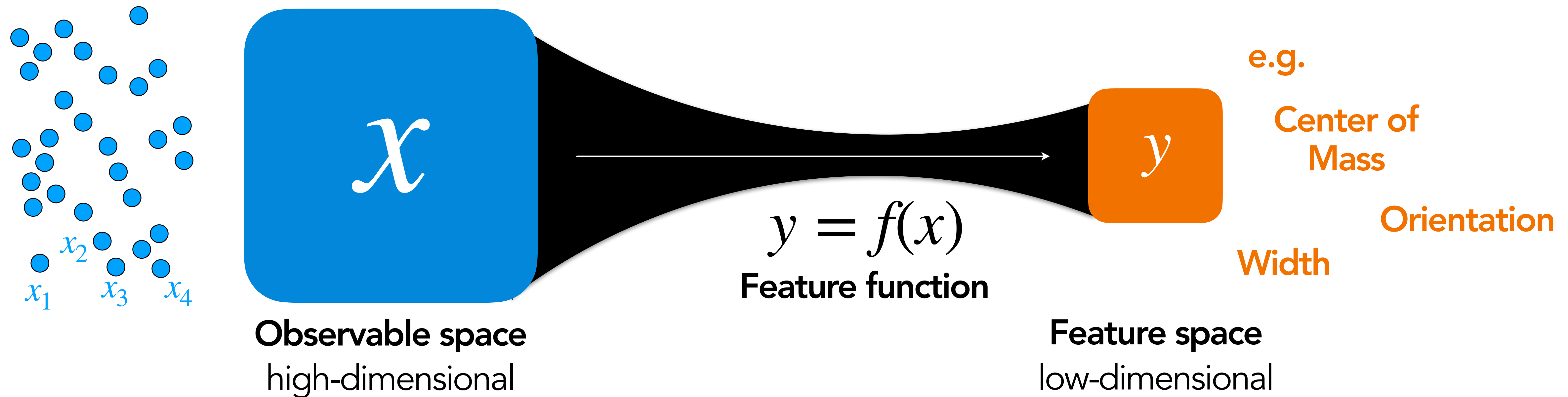
MAX PLANCK INSTITUTE  
for the science of light

# HOW CAN WE DESCRIBE A PHYSICAL SYSTEM WITH ONLY A FEW QUANTITIES?



e. g. Fluid dynamics, Thermodynamics, ...

# FEATURES

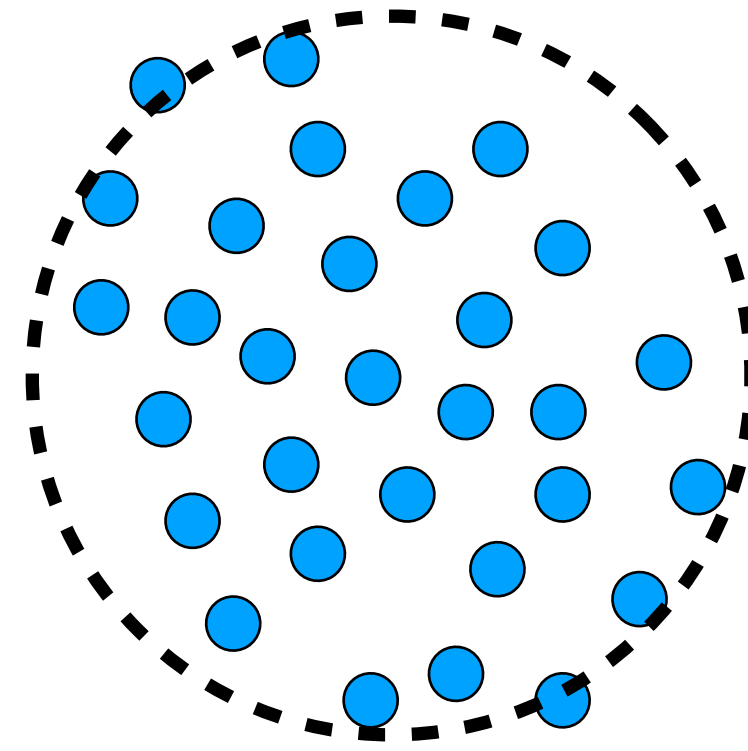
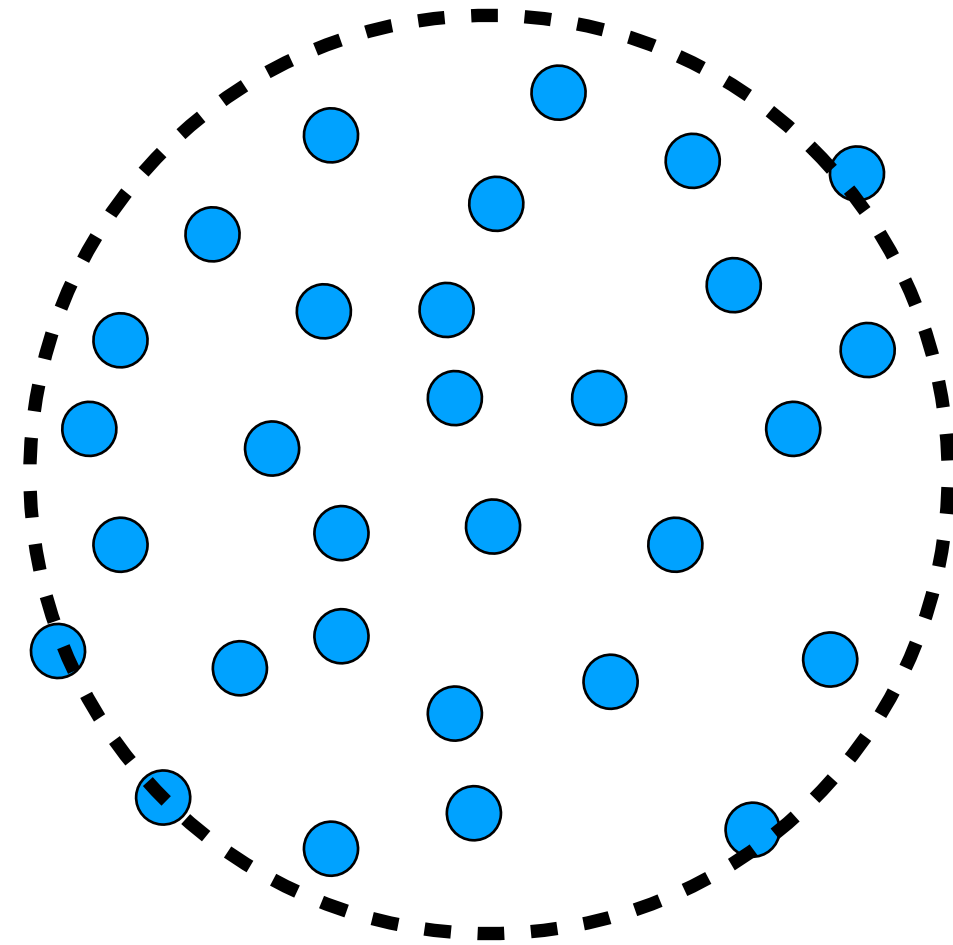
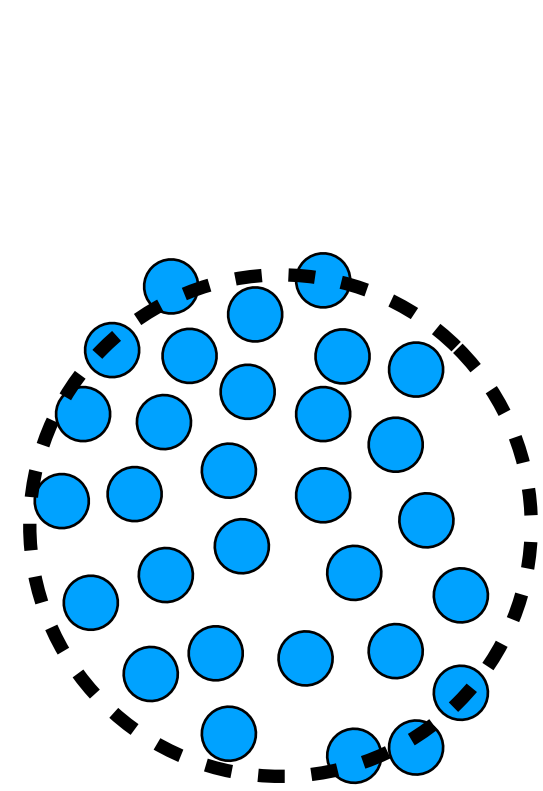


**Goal:** Given a high-dimensional system  $x$ , we want to automatically find the best low-dimensional features  $y = f(x)$  to describe it.

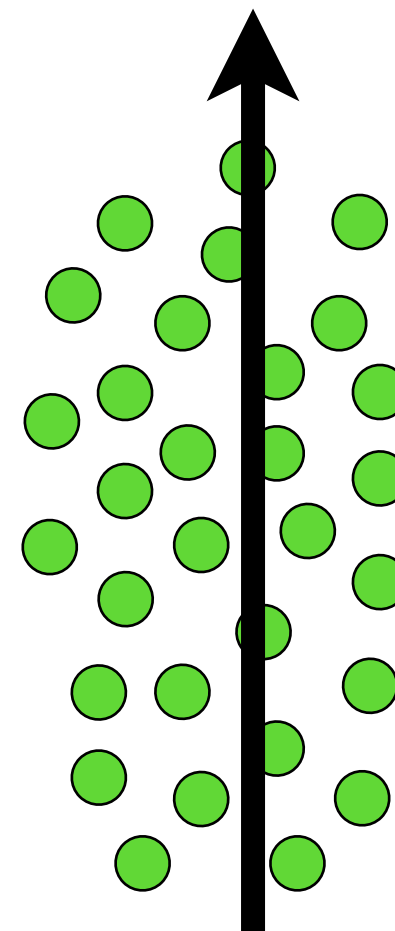
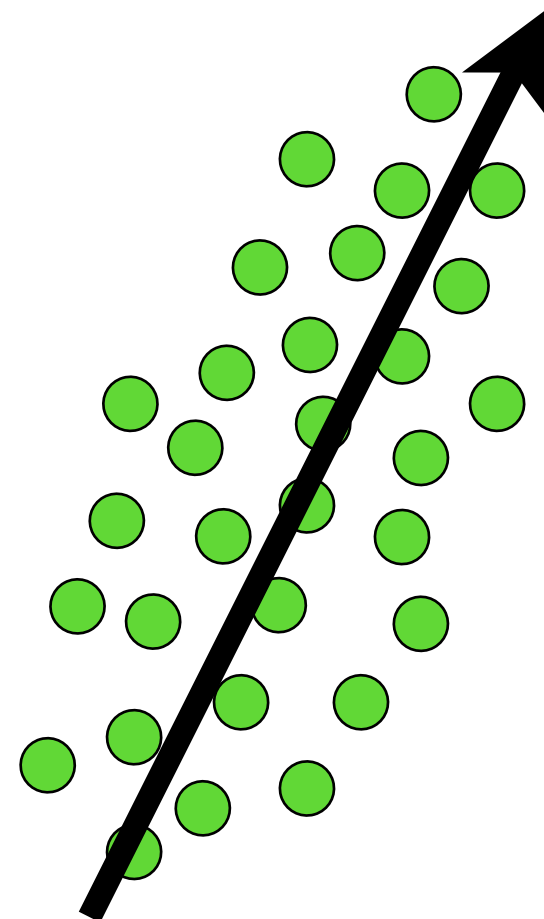
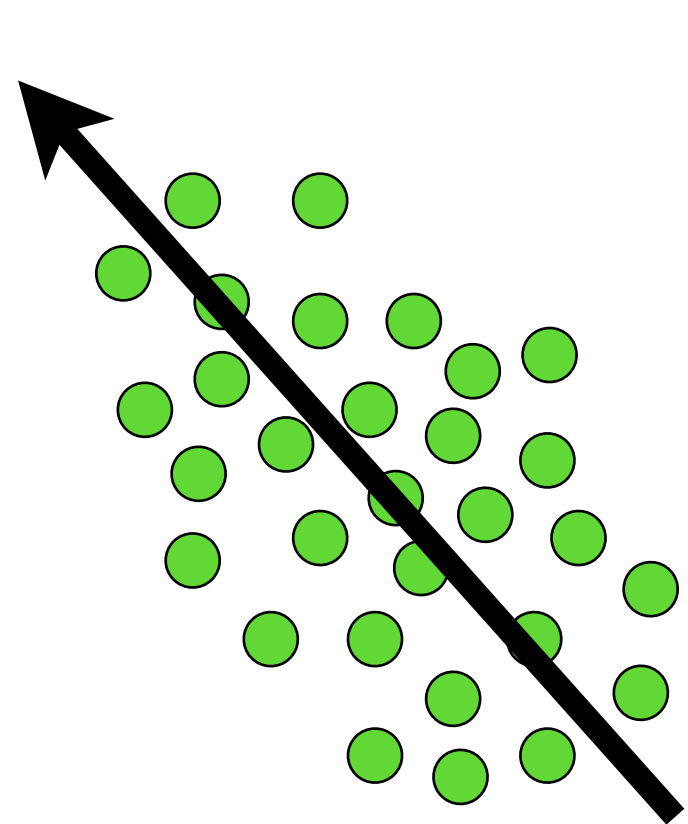


# FEATURE EXTRACTION

Statistical analysis of many observations



Width



Orientation



# MUTUAL INFORMATION

*If two random variables  $x$  and  $y$  are dependent one another, and we are provided with the value of  $y$ , how much do we learn on  $x$ ?*

$$I(x, y) = \overset{\text{Entropy}}{\underset{\text{Conditional Entropy}}{H(y) - H(y|x)}} = \int dx dy P(x, y) \log \frac{P(x, y)}{P_x(x)P_y(y)}$$

- It quantifies the dependence between two random variables
- Always positive, zero iff the variables are independent

**Feature Extraction:** choose  $f(x)$  such that  $I(x, y = f(x))$  is maximized

References:

- *Mutual Information* - e. g. Papoulis, Athanasios, and S. Unnikrishna Pillai. *Probability, Random Variables, and Stochastic Processes*
- Bell, Sejnowsky. "InfoMax: An Information-maximisation Approach to Blind Separation and Blind Deconvolution";

# PROBLEM:

**MUTUAL INFORMATION IS  $+\infty$  FOR ANY FEATURE!**

$$H(y | x) = - \int dx dy P_x(x) P(y | x) \log P(y | x)$$

A continuous deterministic feature  $y = f(x)$  has  $P(y | x) = \delta(y - f(x))$ .

$$H(y = f(x) | x) = - \int dx dy P_x(x) \delta(y - f(x)) \log \delta(y - f(x)) = - \log \delta(0)$$

**it always diverges in this case!**

**$-\infty$**

References:

- Bell, Sejnowsky. "InfoMax: An Information-maximisation Approach to Blind Separation and Blind Deconvolution";
- Gabri  et al. "Entropy and Mutual Information in Models of Deep Neural Networks." *Journal of Statistical Mechanics*;
- Saxe et al. "On the Information Bottleneck Theory of Deep Learning"

**Long-standing problem  
but unsatisfying solutions**

# RENORMALIZED MUTUAL INFORMATION

## our solution

We add Gaussian noise  $\lambda$  to the  $x$  and define a new finite quantity

$$I_\varepsilon(x, y) = I(x, y = f(x + \varepsilon\lambda))$$

We perform the zero-noise limit:

$$\tilde{I}(x, y) = \lim_{\varepsilon \rightarrow 0} I_\varepsilon(x, y) + H(\varepsilon\lambda) = H(y) - \int dx P_x(x) \log \sqrt{|\nabla f(x) \cdot \nabla f(x)|}$$

- Finite quantity and well-defined
- Invariant under feature reparametrization  
(an invertible transformation  $y' = g(y)$  does not change  $\tilde{I}$ )

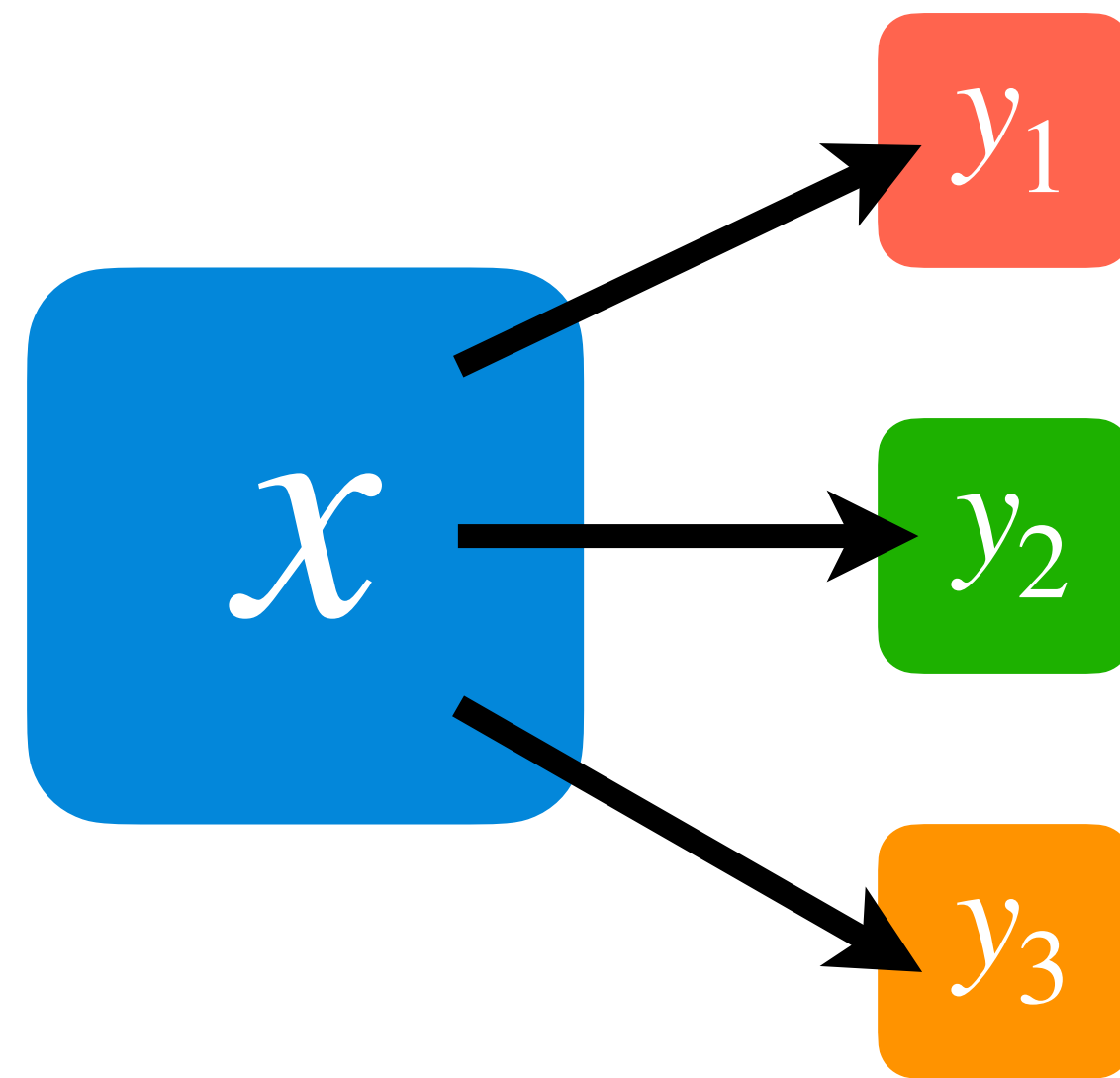
*L. Sarra,  
A. Aiello,  
F. Marquardt*

*arXiv:2005.01912*



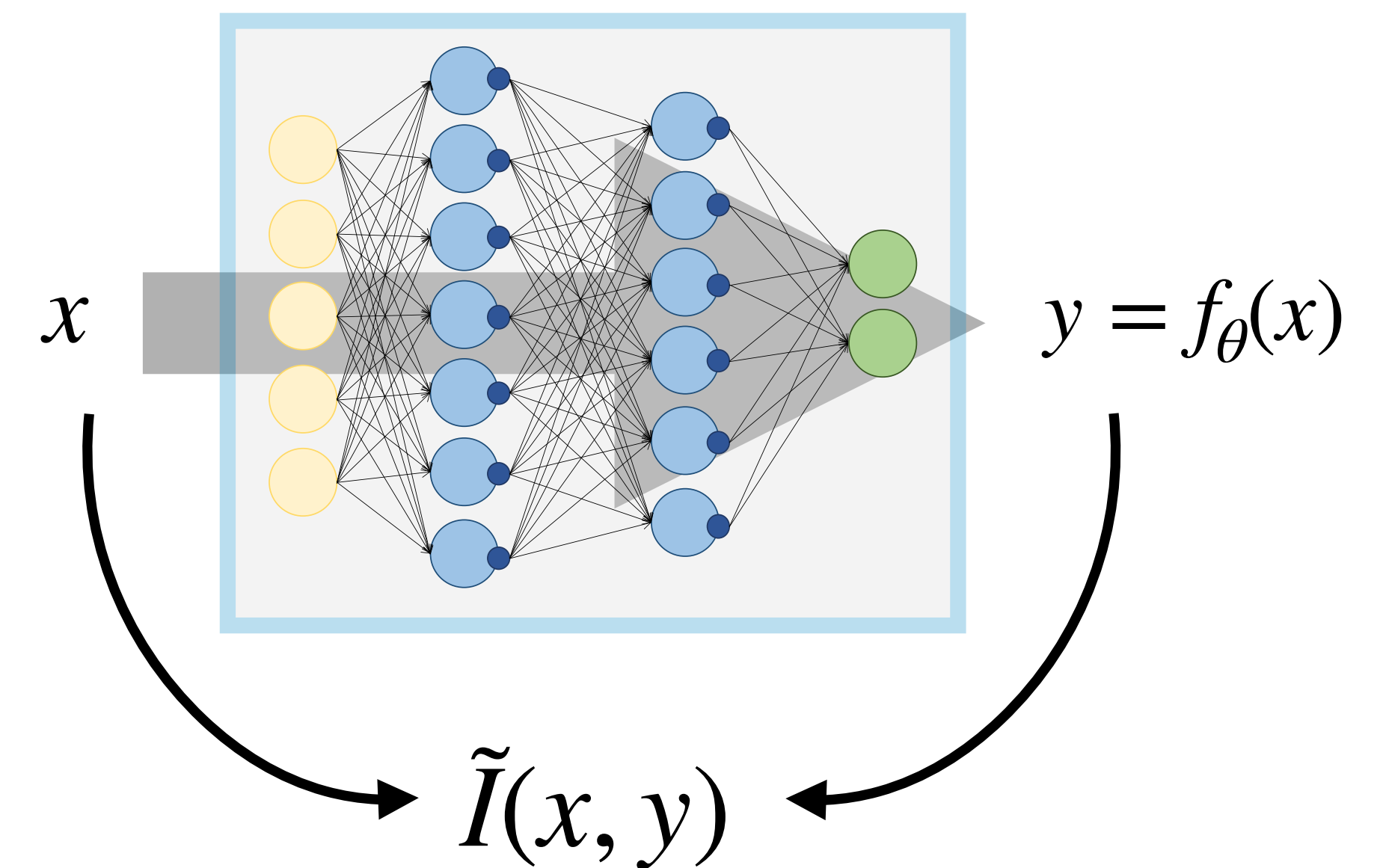
# APPLICATIONS

## Feature Selection



Find out how useful a given macroscopic quantity is to describe the system

## Feature Extraction



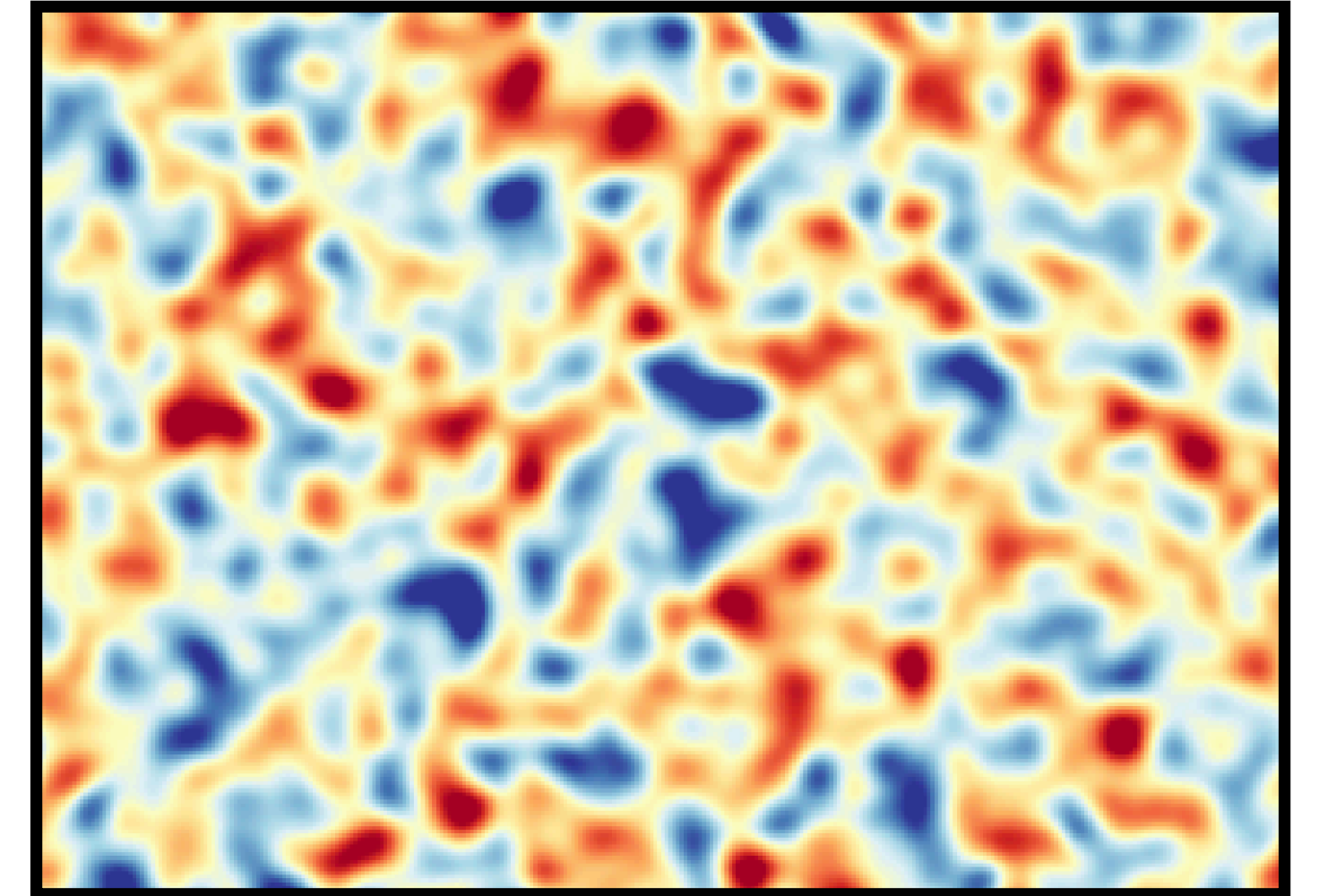
Optimize over a class of functions parametrized by a neural network



# POSSIBLE APPLICATIONS

- Characterize different phases of matter
- Partial prediction of time evolution
- Discover residual regularities of chaotic systems
- Study of intermediate layers of neural networks during training
- Unsupervised representation learning
- ...

Fluctuating fields



Many-particle systems

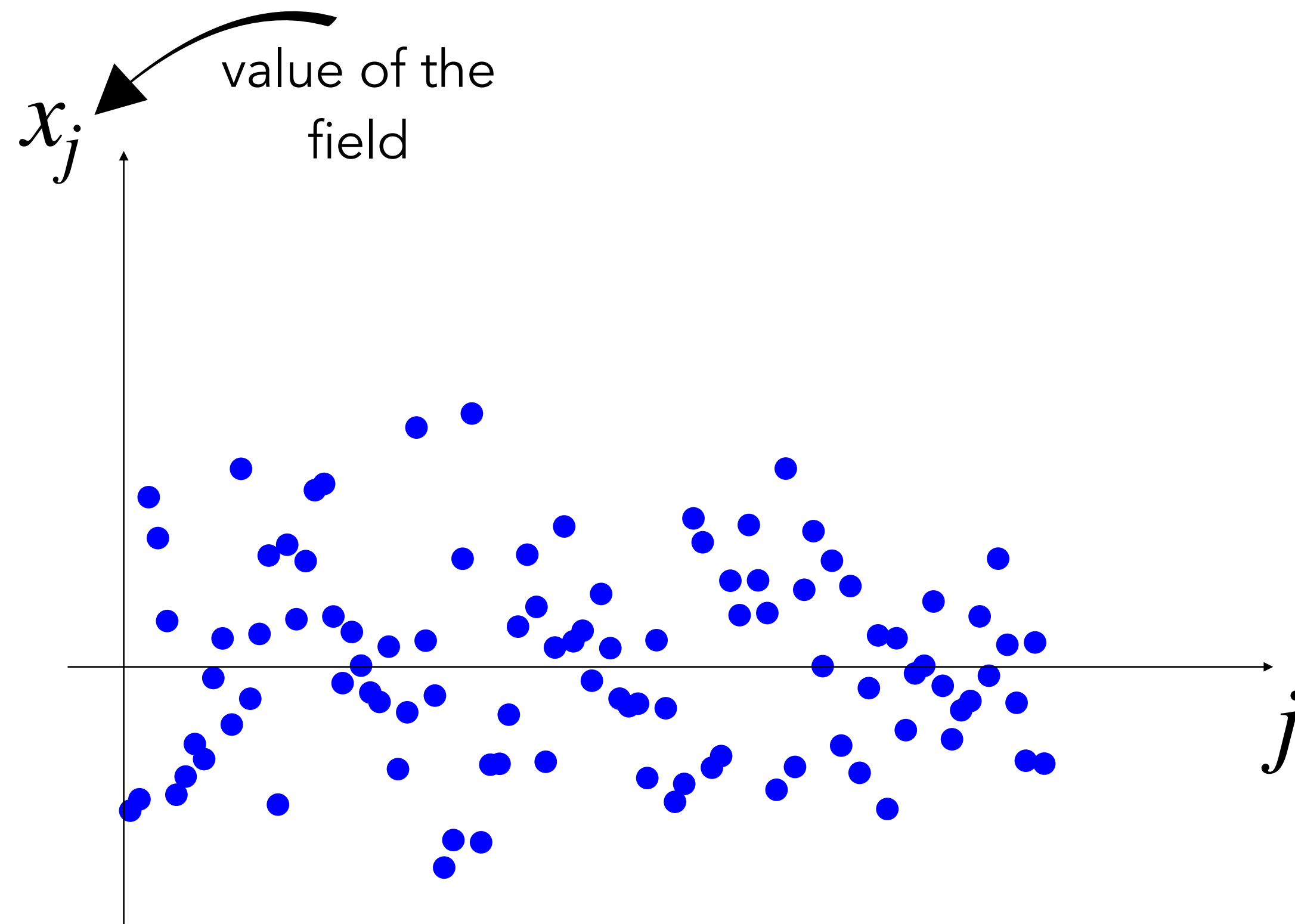


ESA/Hubble

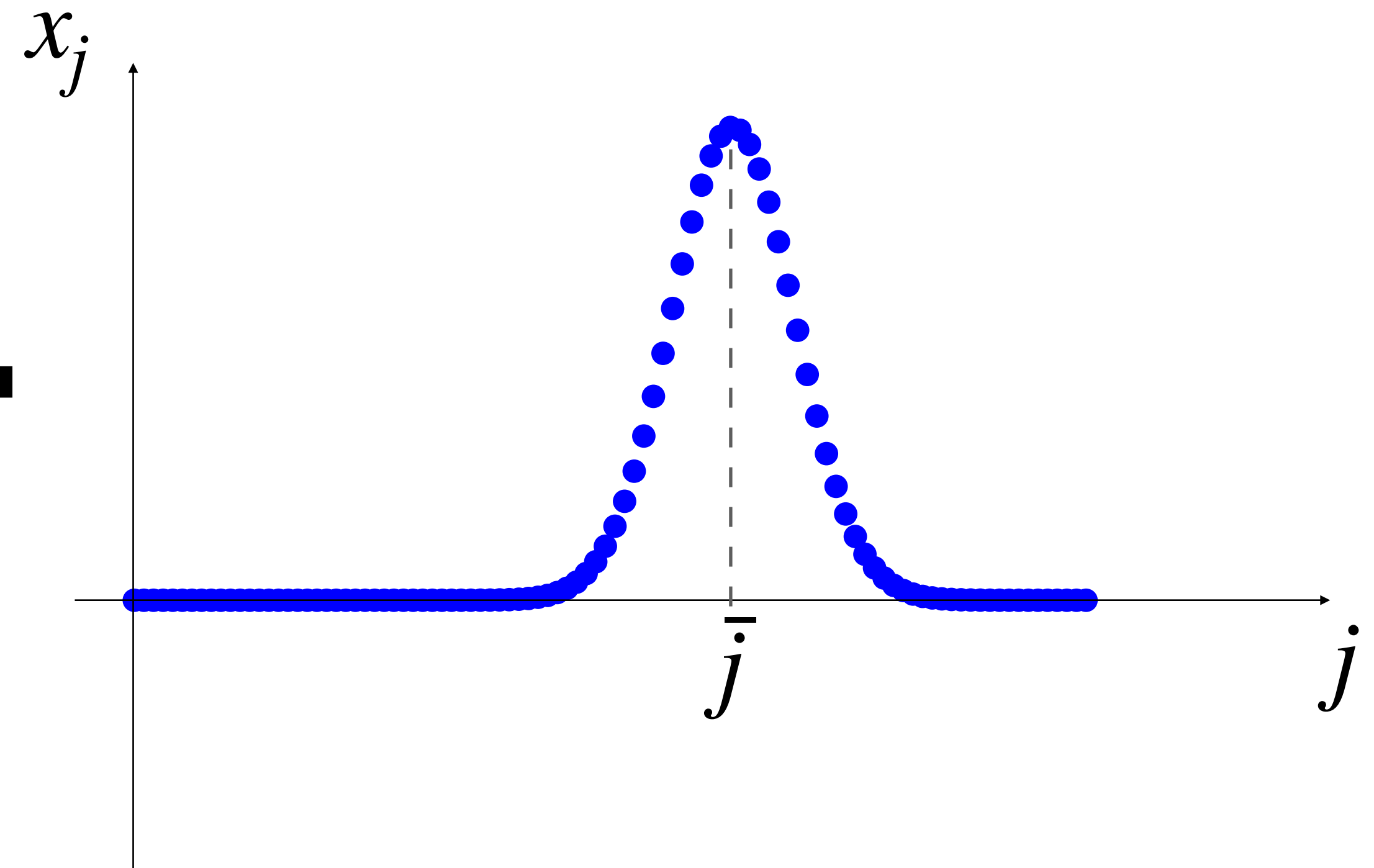


# EXAMPLE

## Fluctuating field with wave packet



+



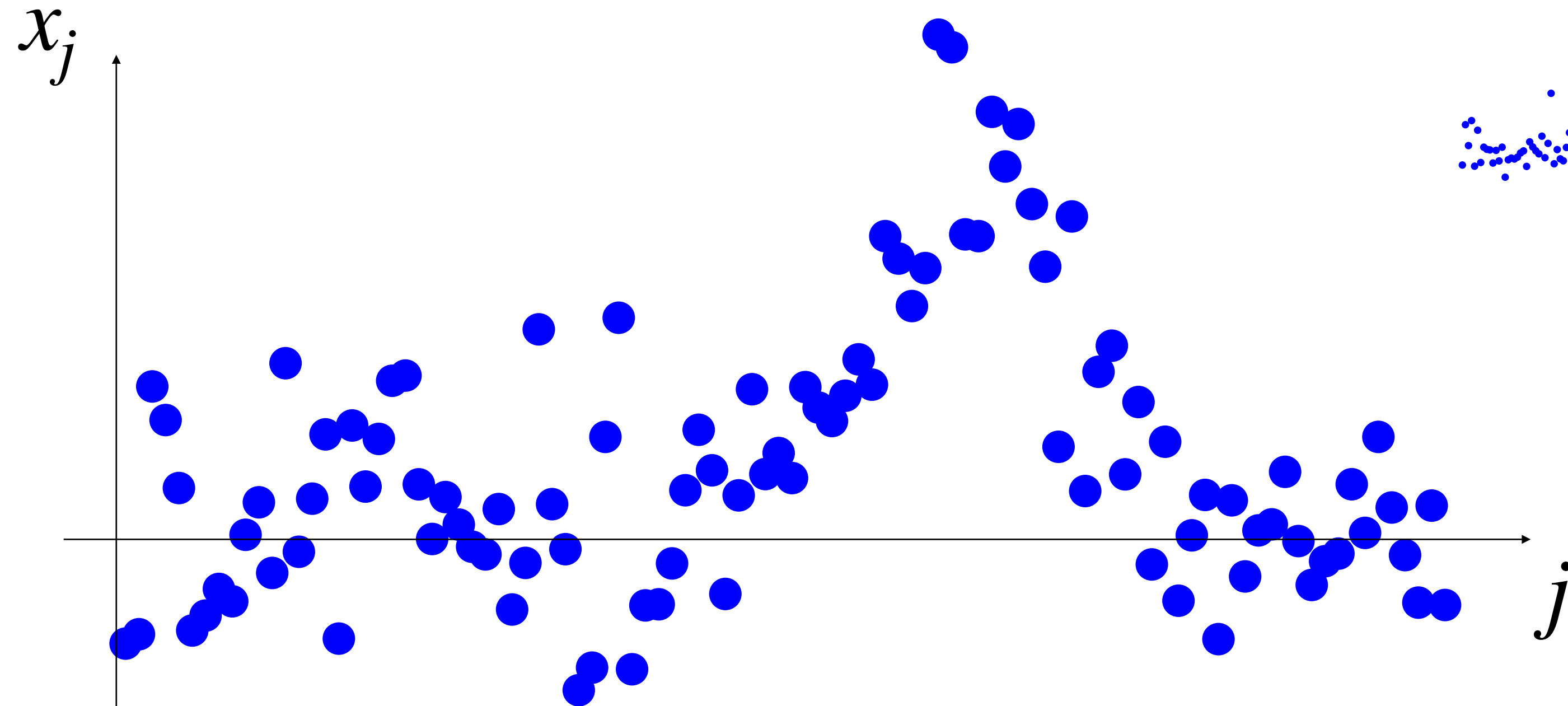
thermal noise

wave packet( $j - \bar{j}$ )

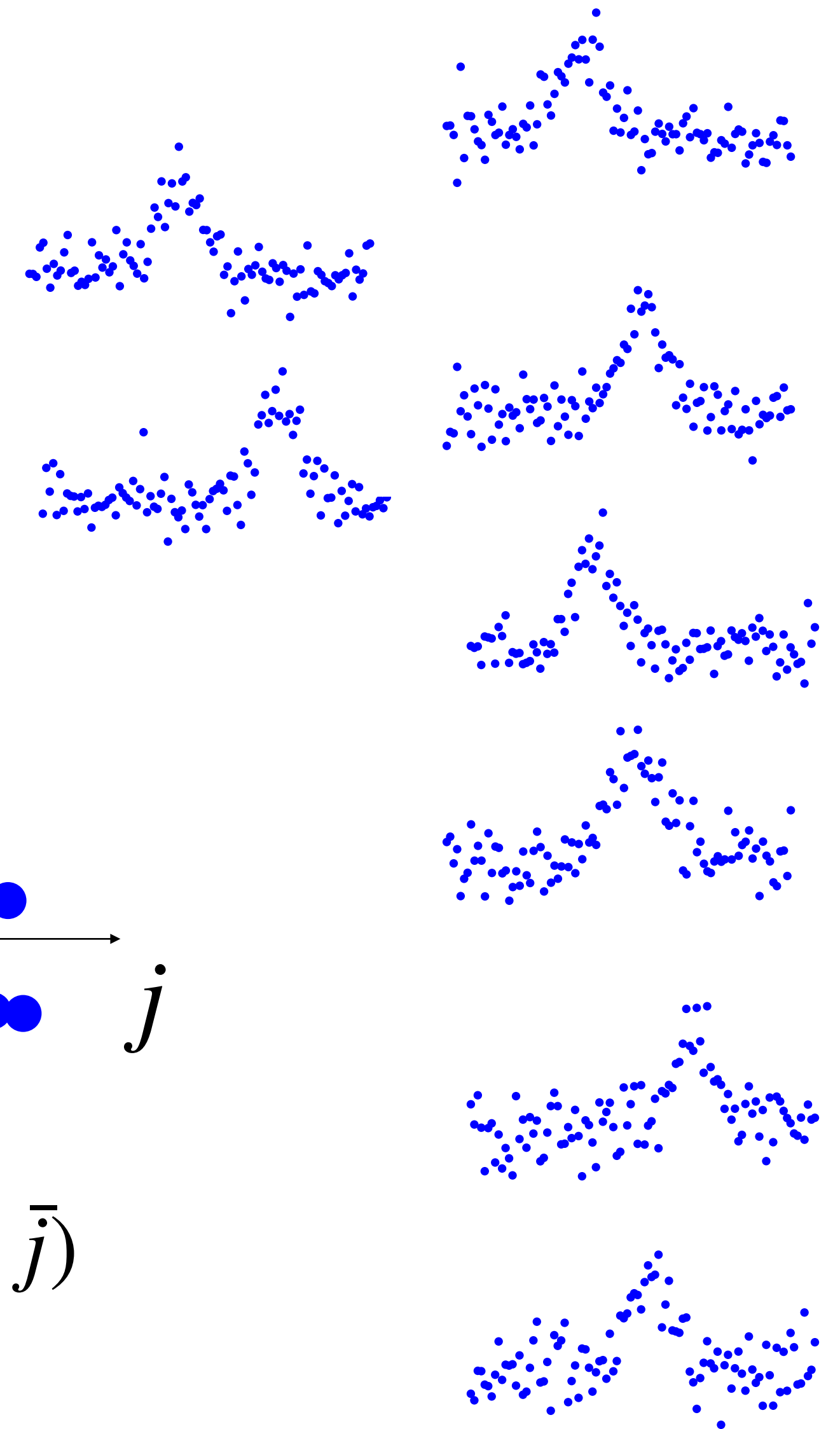


# EXAMPLE

## Fluctuating field with wave packet



$$x_j = \text{thermal noise} + \text{wave packet}(j - \bar{j})$$



# EXAMPLE

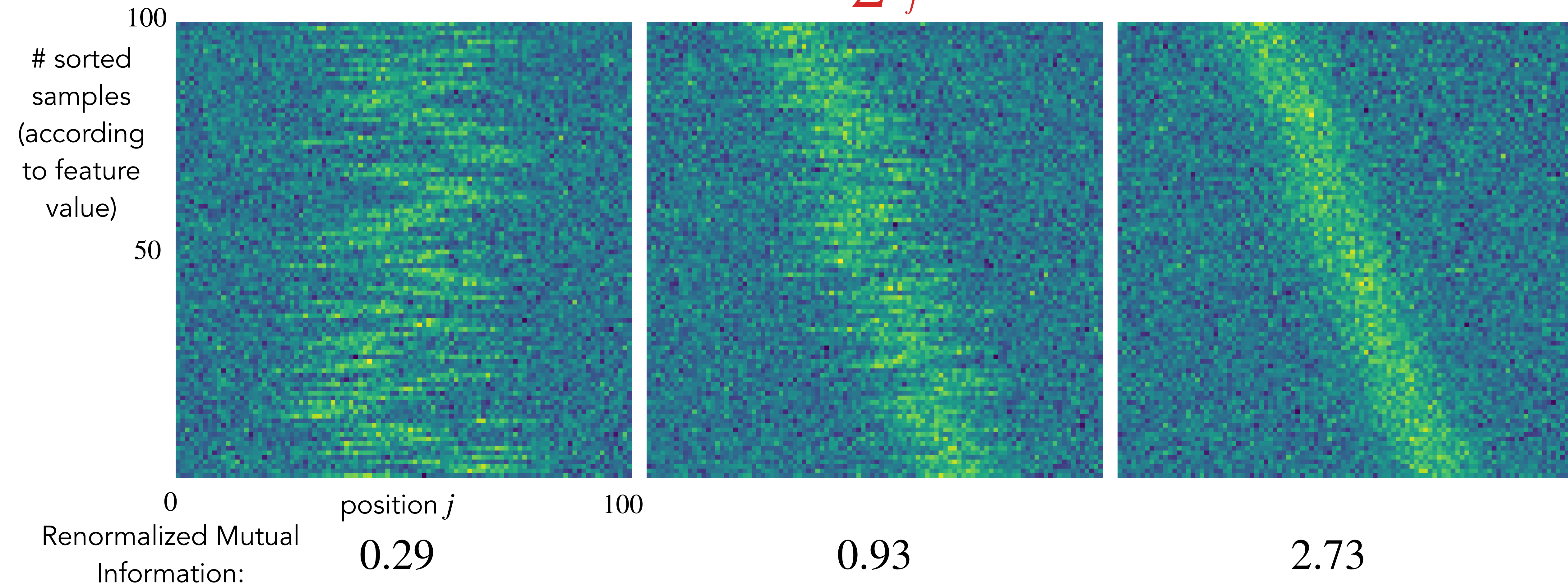
## Fluctuating field with wave packet

$x_j$

$$y = \sum x_j^2$$

$$y = \frac{\sum jx_j^2}{\sum x_{j'}^2}$$

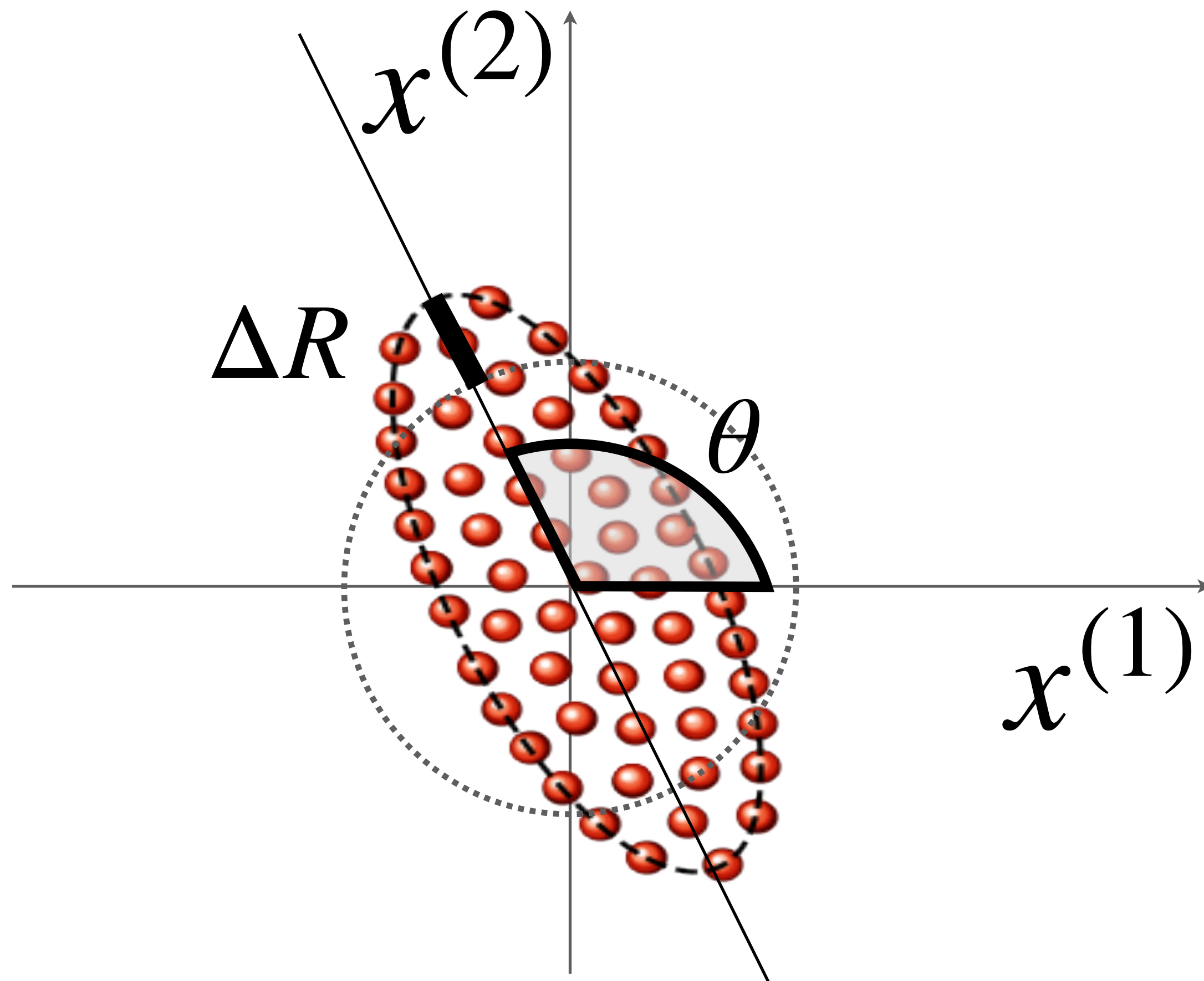
$$y = f_{NN}(x)$$



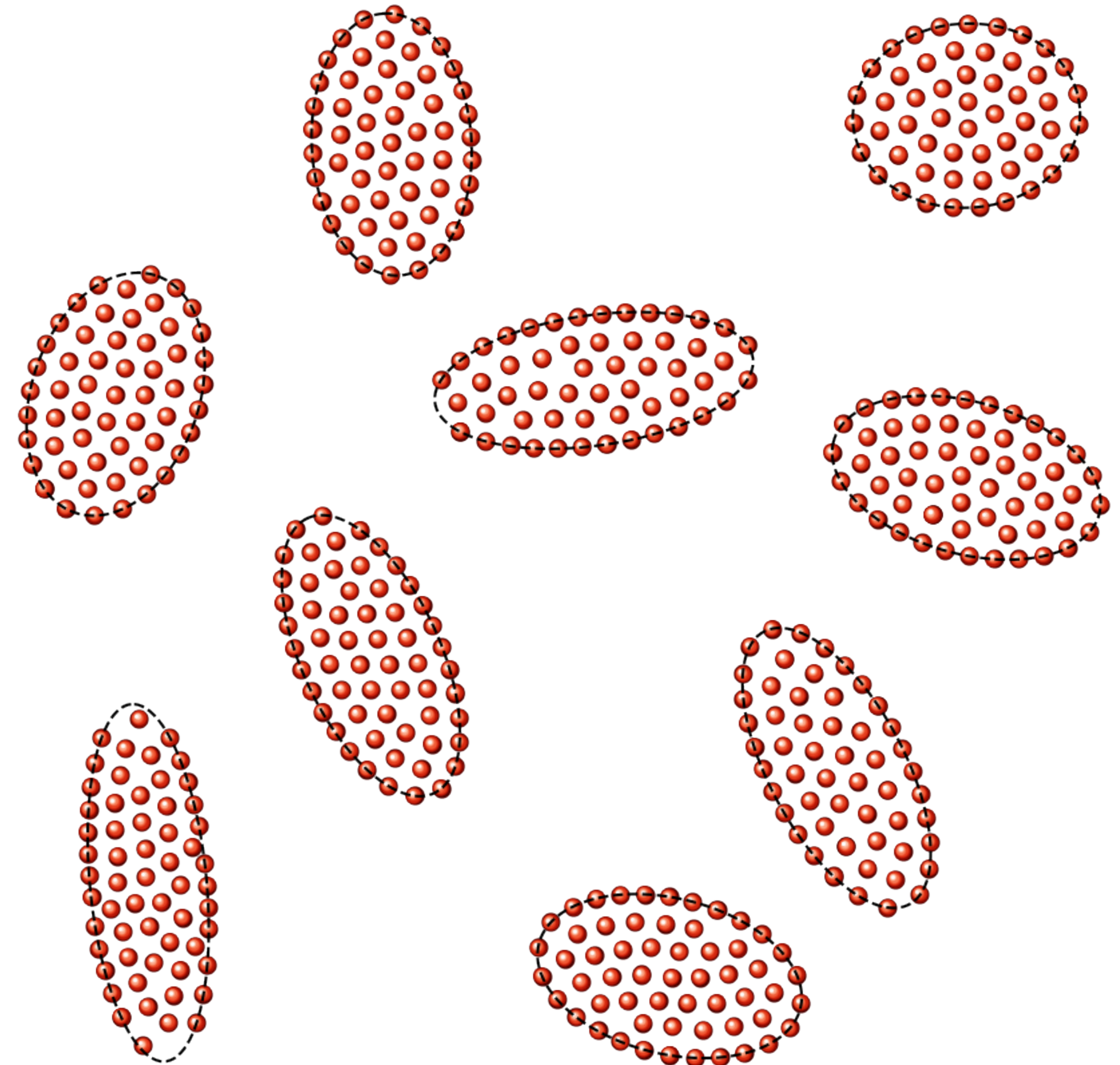
Reference: in our paper *arXiv:2005.01912*, we also estimate the quality of the feature representations in a more quantitative way by comparing the performance on a supervised regression task

# EXAMPLE

## Liquid Drop



Each ellipse has a different  
**deformation** and **orientation**

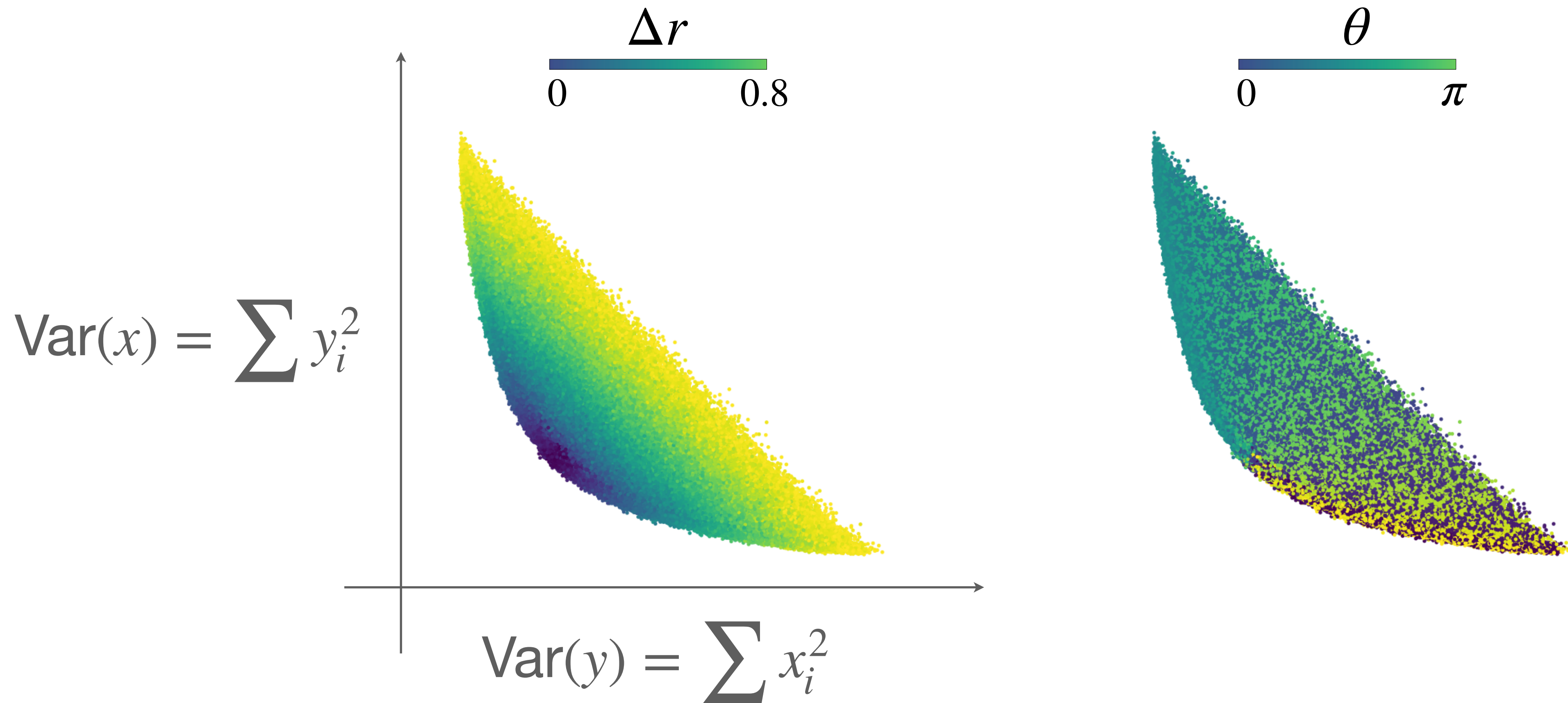




# EXAMPLE

## Liquid Drop

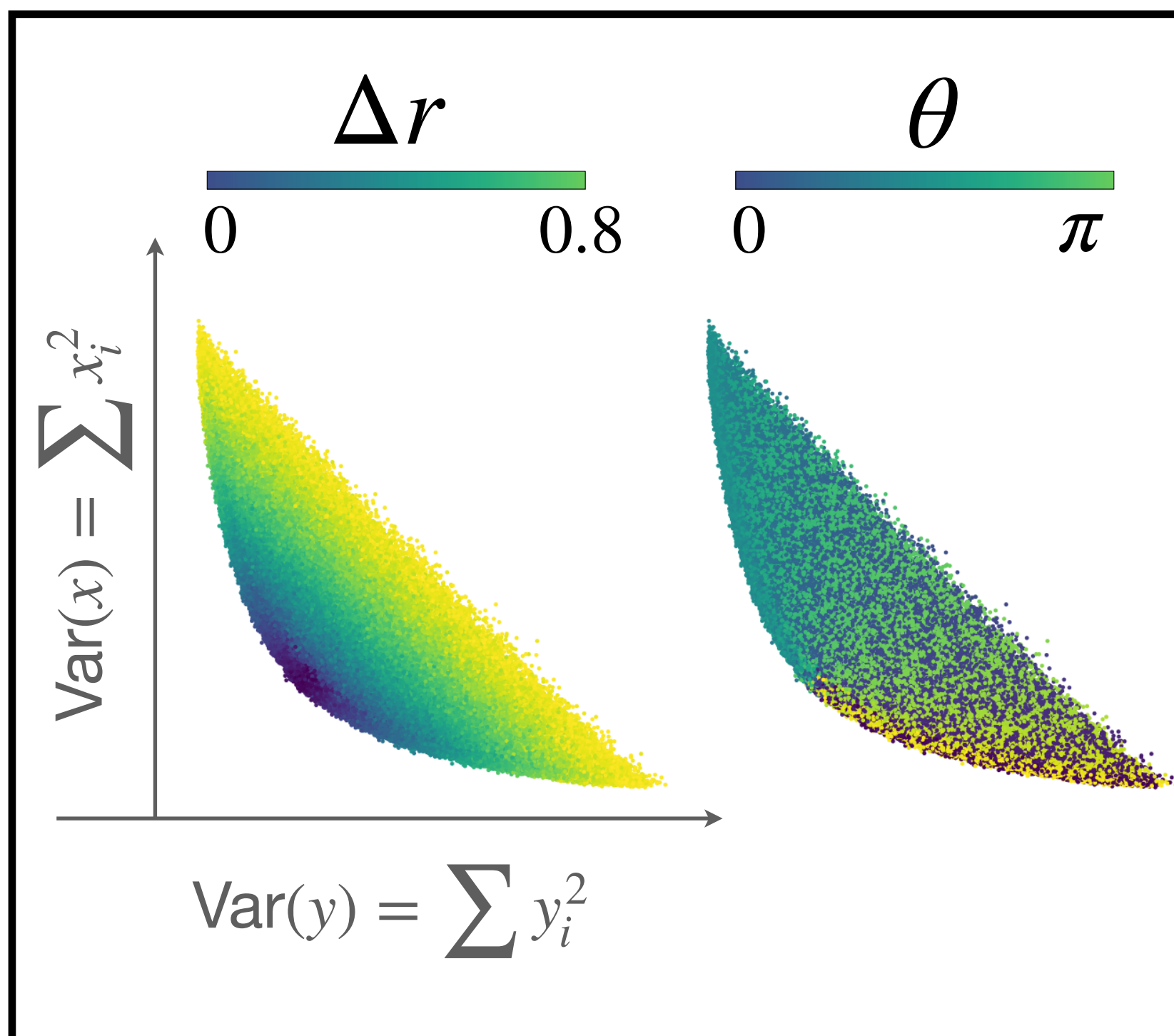
Handcrafted Features



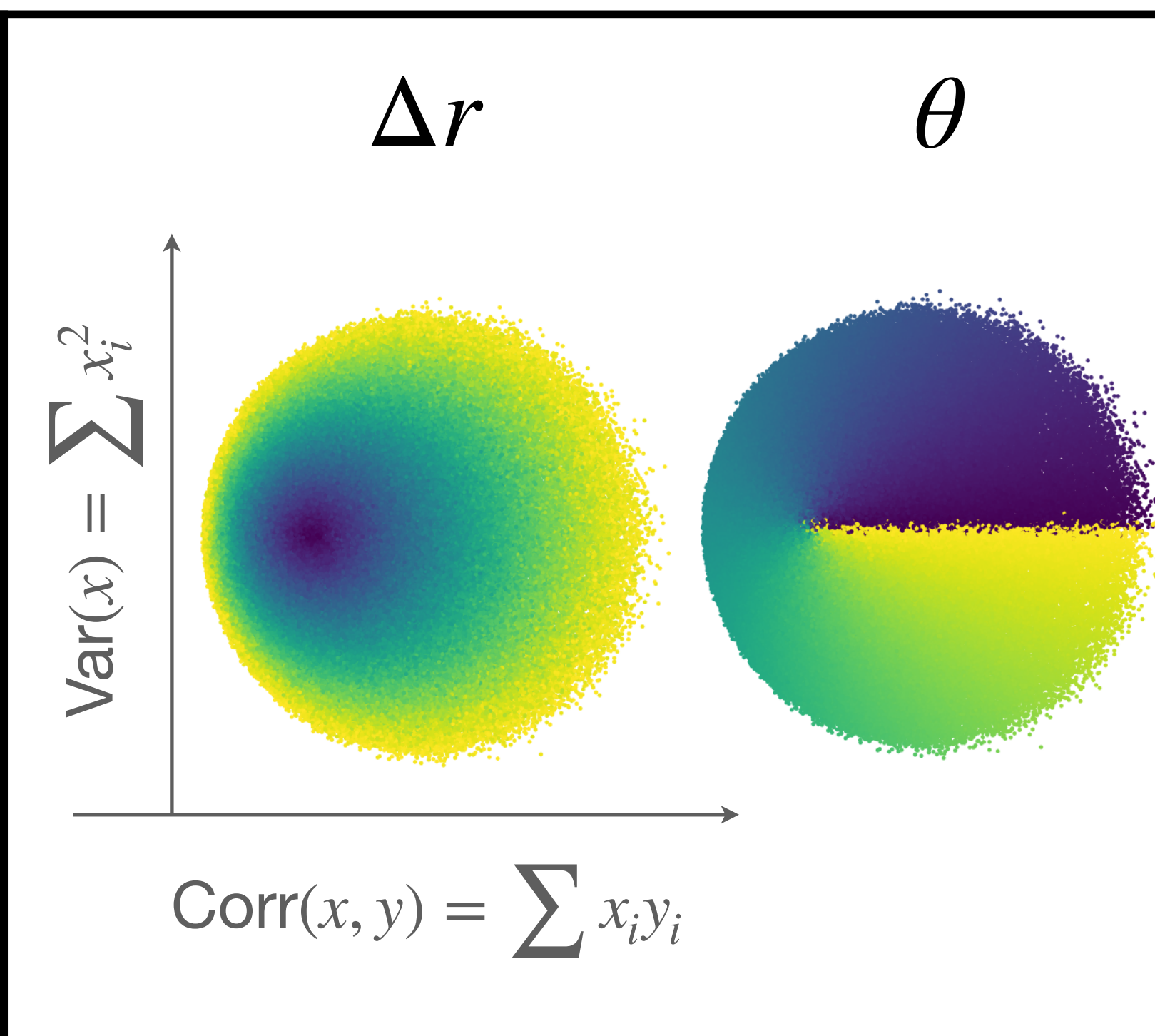
# EXAMPLE

## Liquid Drop

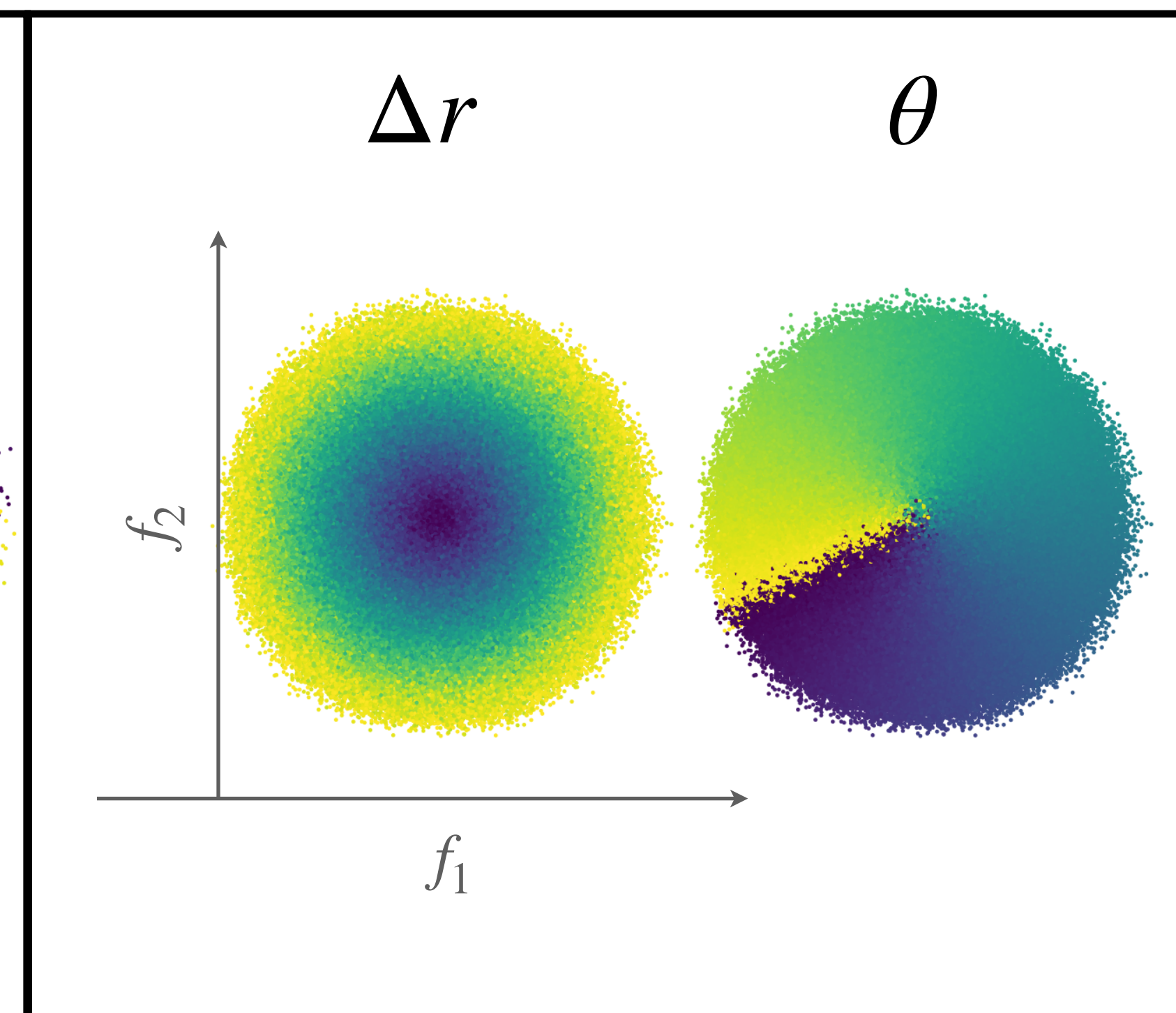
Handcrafted Feature



Handcrafted Feature



Extracted Feature (neural network)



Renormalized Mutual Information:

1.75

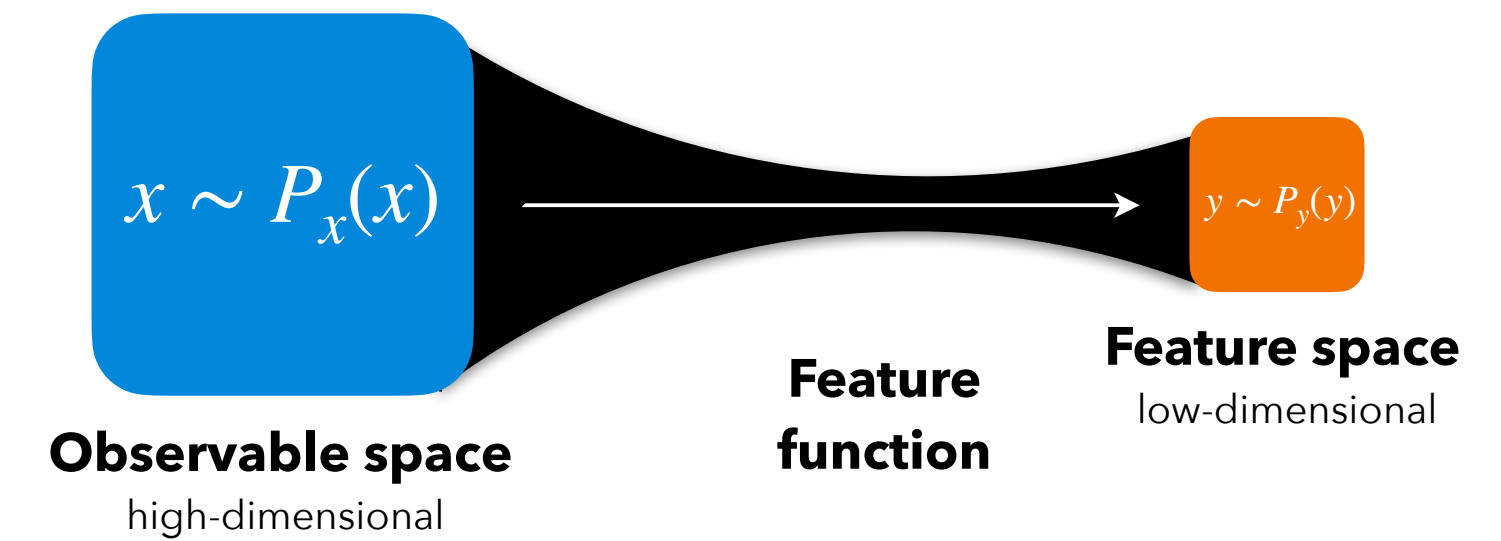
3.04

3.22

Reference: in our paper [arXiv:2005.01912](#), we also estimate the quality of the feature representations in a more quantitative way by comparing the performance on a supervised regression task



# OUTLOOK



What are the most relevant features of a system?

$$\tilde{I}(x, y) = H(y) - \int dx P_x(x) \log \sqrt{|\nabla f(x) \cdot \nabla f(x)|}$$

Many possible applications in Physics and Machine Learning

- Characterize different phases of matter
- Partial prediction of time evolution
- Discover residual regularities of chaotic systems
- Study of intermediate layers of neural networks during training
- Unsupervised representation learning
- ...



"Renormalized Mutual Information for Artificial Scientific Discovery",  
LS, Andrea Aiello, and Florian Marquardt,  
Phys. Rev. Lett. **126**, 200601