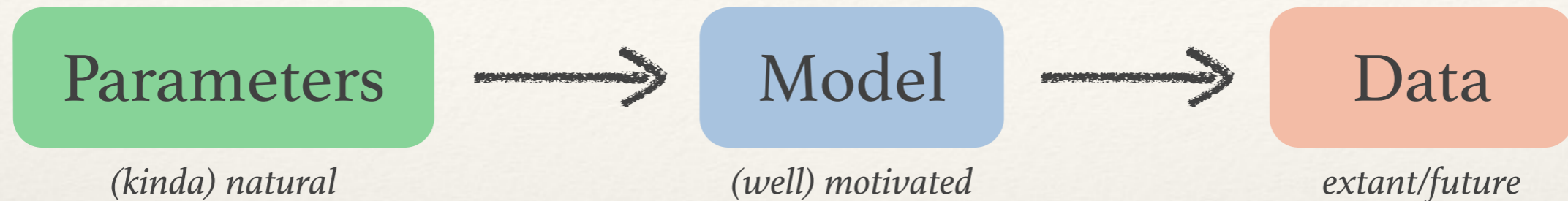

Machine Learning
Beyond Physics
from a Physicist's Perspective

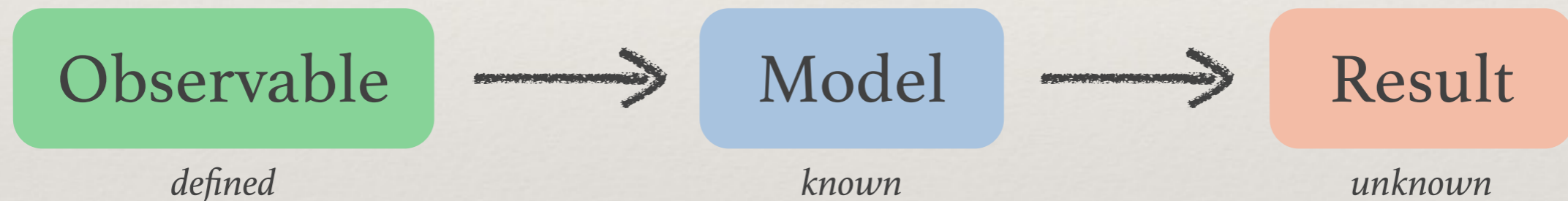
Jonathan Walsh

How do we do physics?

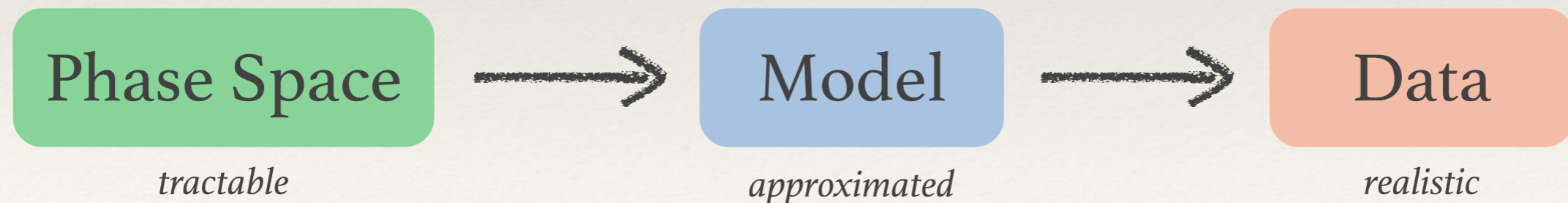
model building:



calculations:



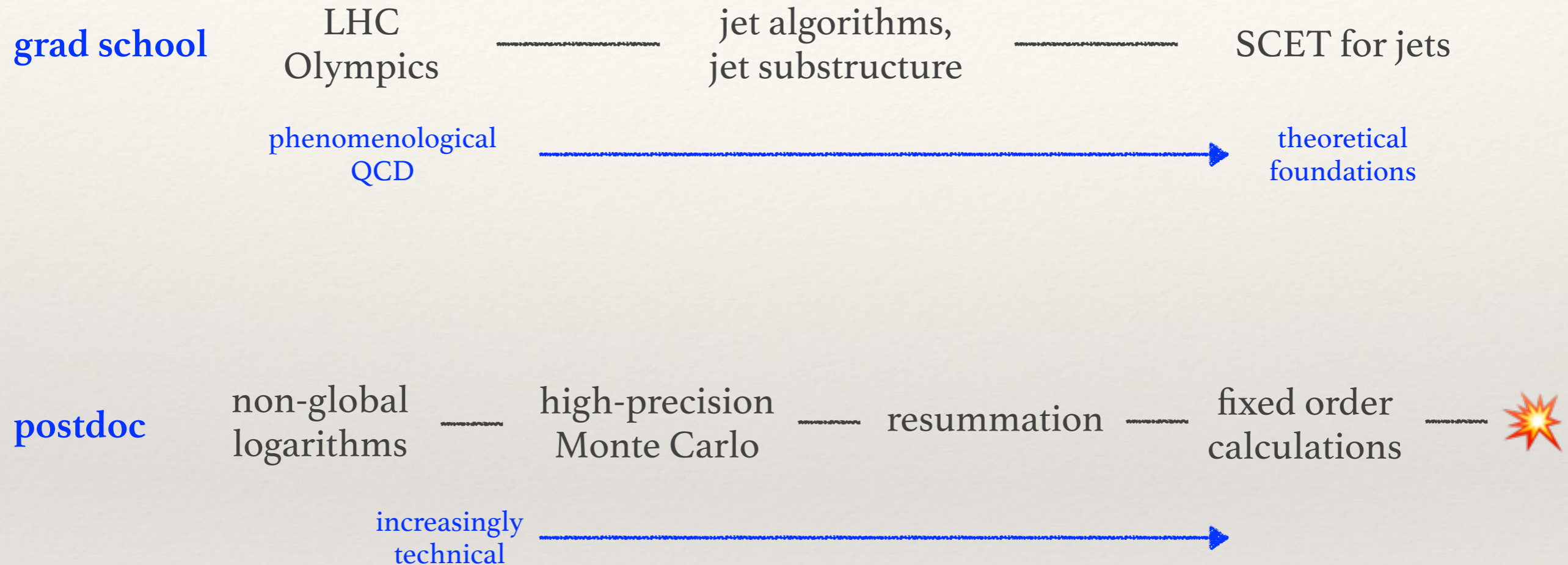
simulations:



A core tenet of particle physics

We can
understand
the model

My own experience in physics (QCD-side)



Evolution of my research program largely governed by the process:

1. finding interesting problems in QCD
2. **understanding the underlying theory**
3. applying it to more interesting problems
4. rinse and repeat

Modeling beyond physics

Beyond physics, the model may be too complex or just unknown.

A first-principles approach *does not apply* in these cases.

We need tools capable of model inference that can learn and utilize relevant information in the data

2 common approaches: *statistical* and *deterministic* modeling

- statistical : derive the form of the model
- deterministic : input the form of the model

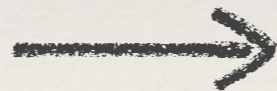
Reframing the core tenet

Beyond physics, the model may be too complex or just unknown.

A first-principles approach *does not apply* in these cases.

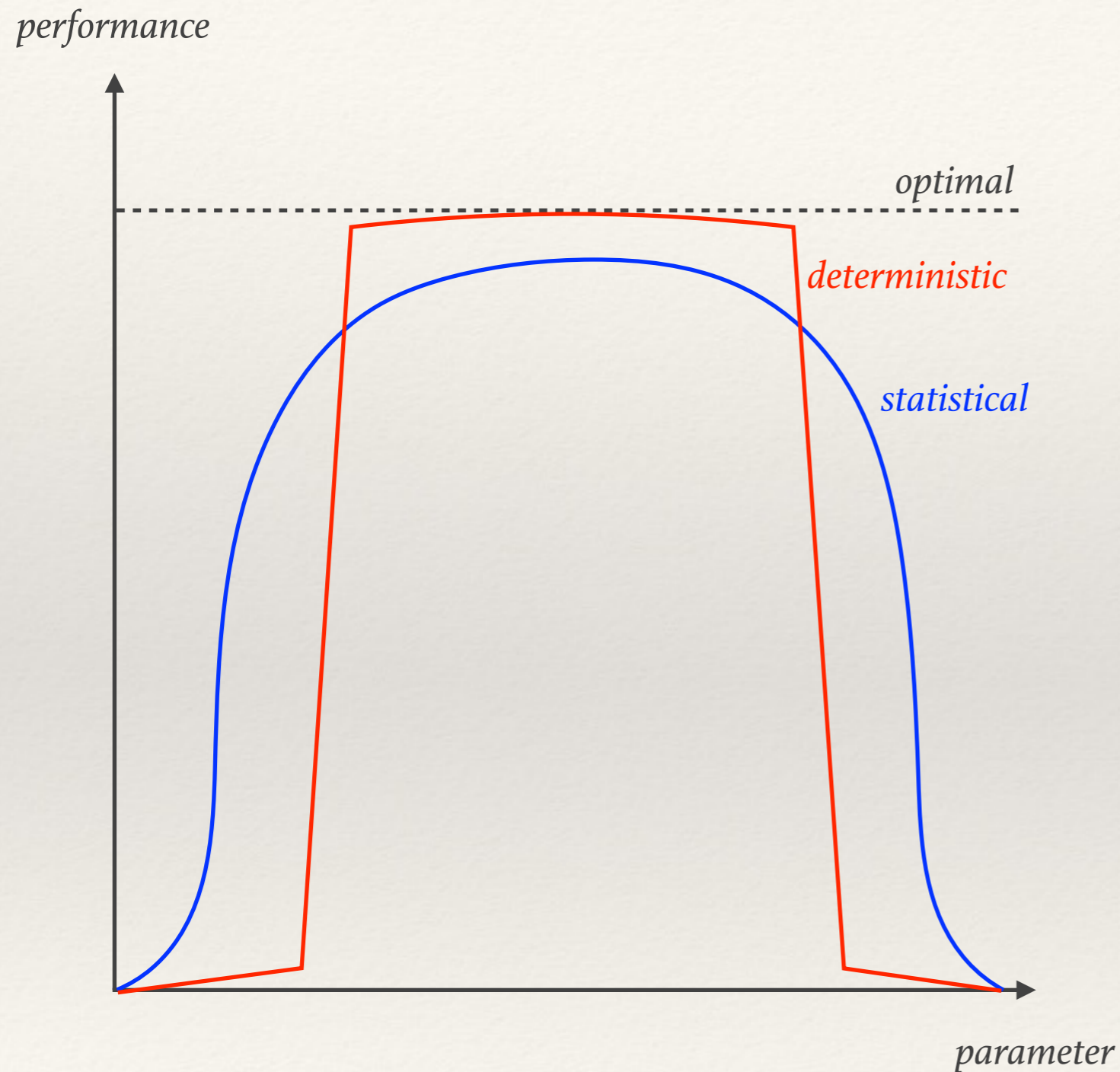
We need tools capable of model inference that can learn and utilize relevant information in the data

**We can
understand
the model**



**We can
understand
the *data***

Statistical vs. Deterministic Modeling



deterministic models:

- high accuracy in a region of phase space
- bad failure modes outside domain of applicability

statistical models:

- reasonable accuracy across phase space
- graceful failure modes

which should we prefer?

Minimax Optimization

optimize for the best performance of the worst case

In most applications, we care about performance in the worst case:

- you want your bike/car/train/plane not to crash
- you want to avoid serious illness
- you buy insurance for costly rare events
- you prioritize products working over their features

→ User experience is most sensitive to performance in the worst cases

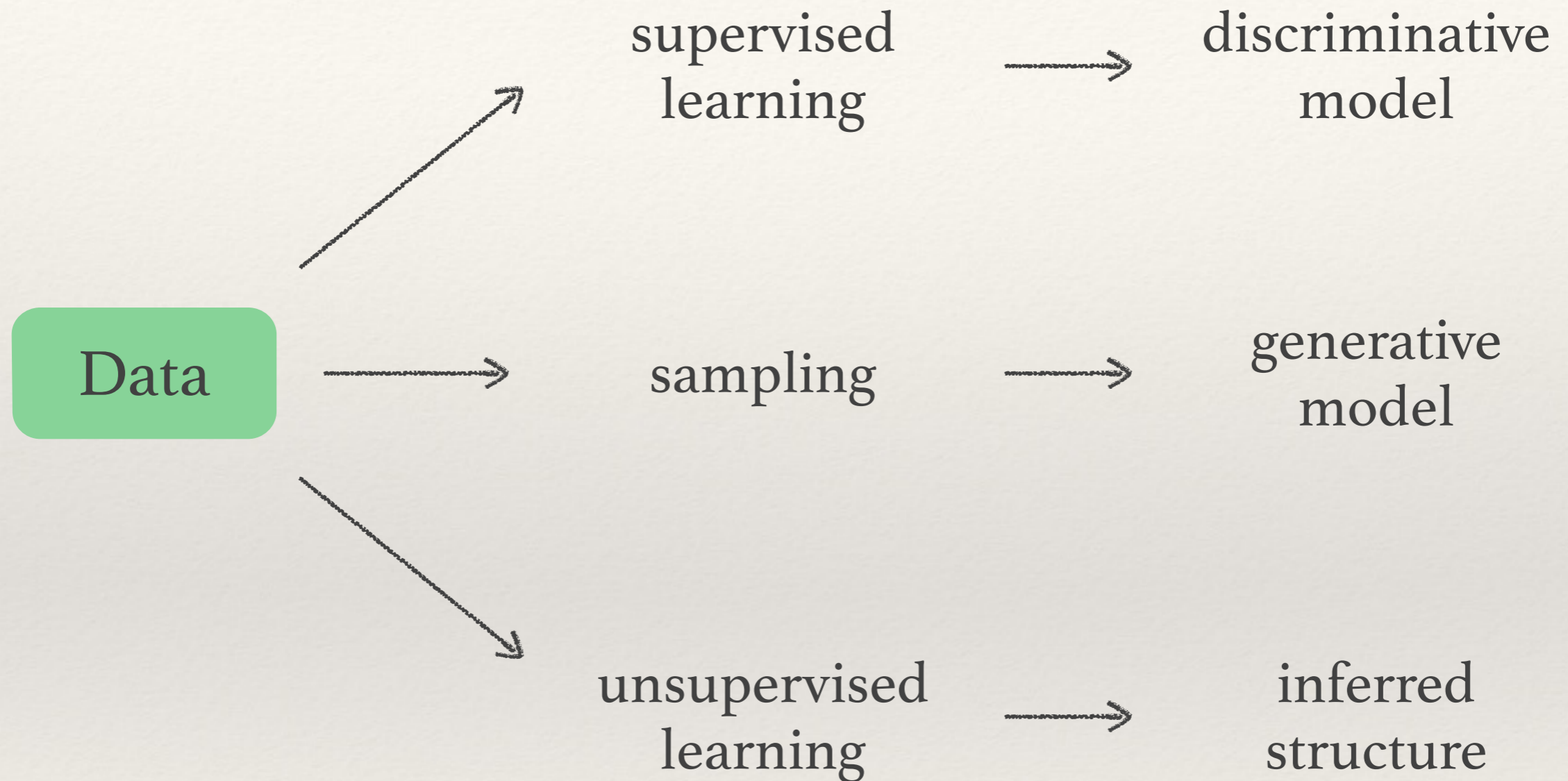
People tend to remember the worst parts of an experience and base a valuation more heavily on that:

- the worst dishes at a meal
- reliability of a car
- annoyances in computer UX (e.g. Mac vs. Windows)
- *everything* about flying

A core tenet of machine learning

Learn expressive models

Understanding Datasets



often, we want to transform the data into features as a first step

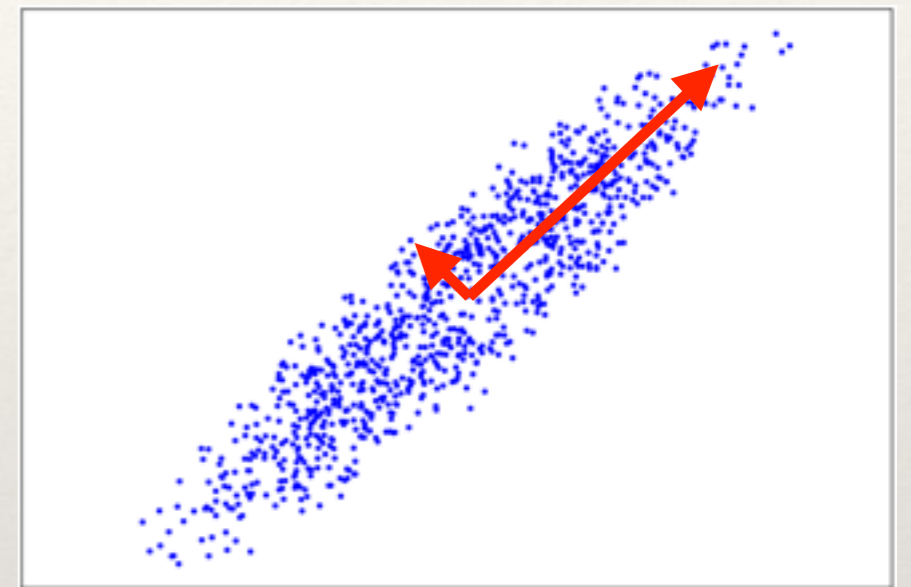
Feature Extraction

feature extraction: start with a linear model

PCA: rotate to a basis which maximizes the variance along principal directions

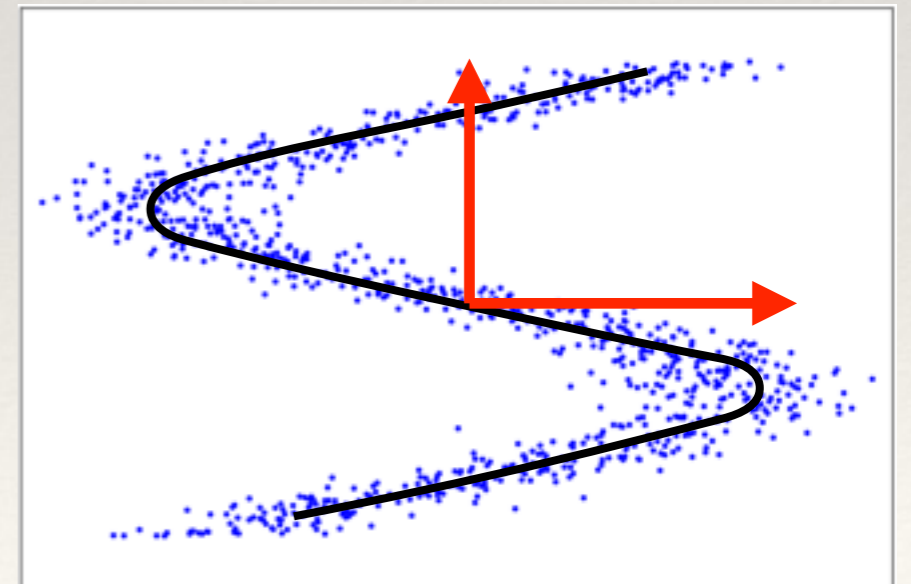
$$v = \mathbf{W}(x - \mu)$$

orthogonal transform sample feature means



can be ineffective for nonlinear manifolds

manifold learning tools:
manifold learning (e.g. isomap, LLE),
autoencoders

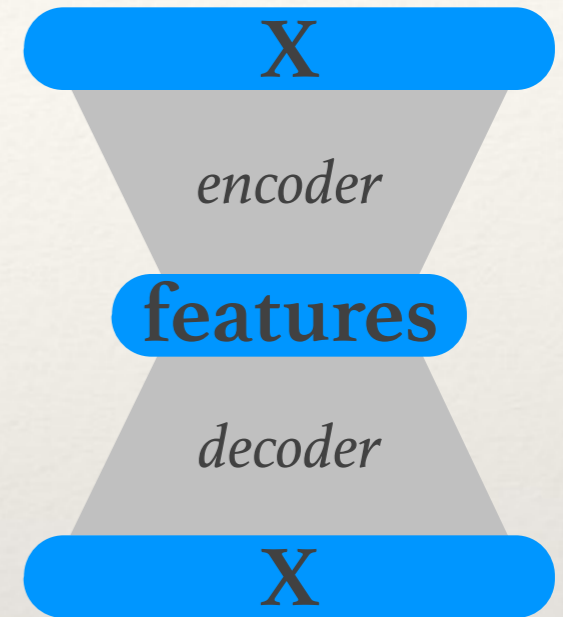


Neural Networks and Autoencoders

autoencoder

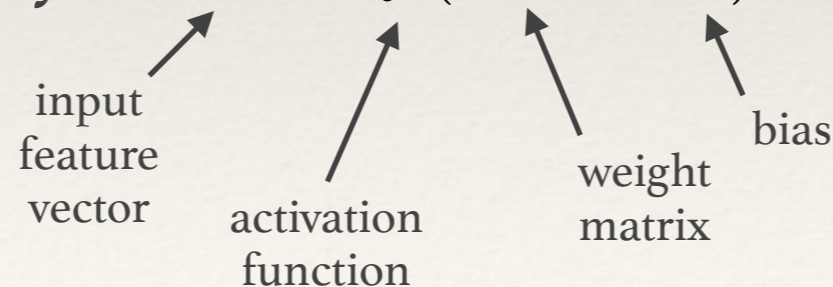
$$\mathcal{M} : X \rightarrow X$$

a model that can reconstruct its inputs
with a constraint on intermediate features
(the encoded representation)



encoder and decoder: (stacks of) neural network layers

fully connected layer: $x \rightarrow f(\mathbf{W}x + b)$



Autoencoders

simplest autoencoder

encoder and decoder: single layers, tied weights

encoder: $y = f(\mathbf{W}x + b)$

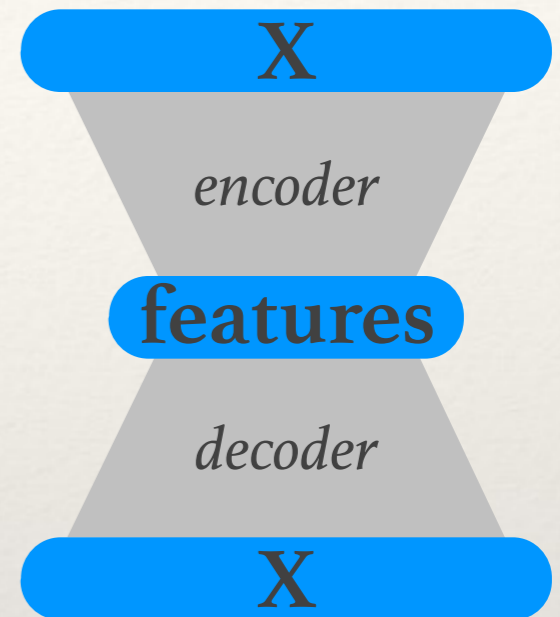
full network:

$$x \rightarrow \mathbf{W}^T y - b = \mathbf{W}^T f(\mathbf{W}x + b) - b$$

network is trained to minimize reconstruction error

with a linear activation, the optimal solution is PCA
(where the bias removes the sample mean)

autoencoders are powerful tools for nonlinear manifold learning

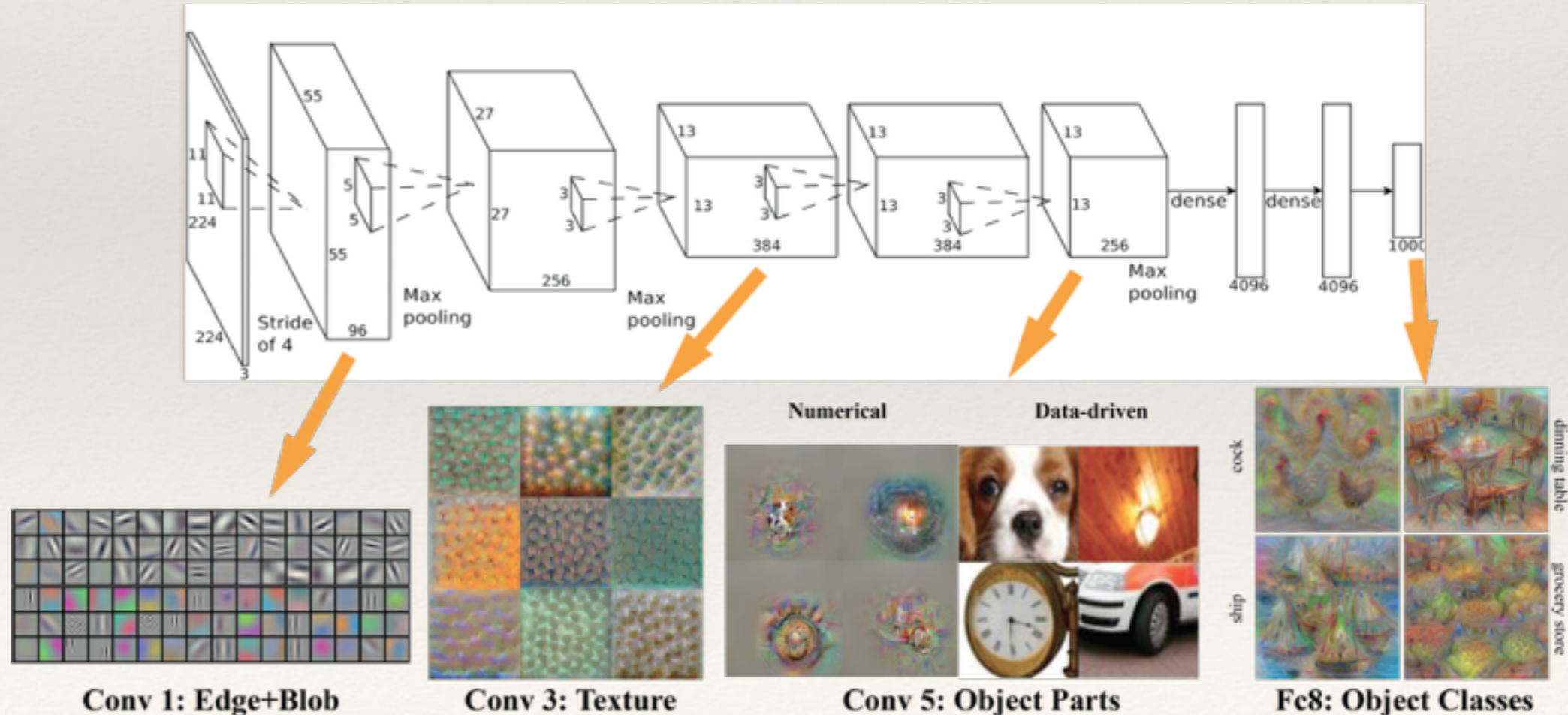


Building Expressive Models

deep learning

models constructed from many simple transformation layers

challenges: effective learning algorithms and architectures, intelligent uses of data



How do we do machine learning?

discriminative models:



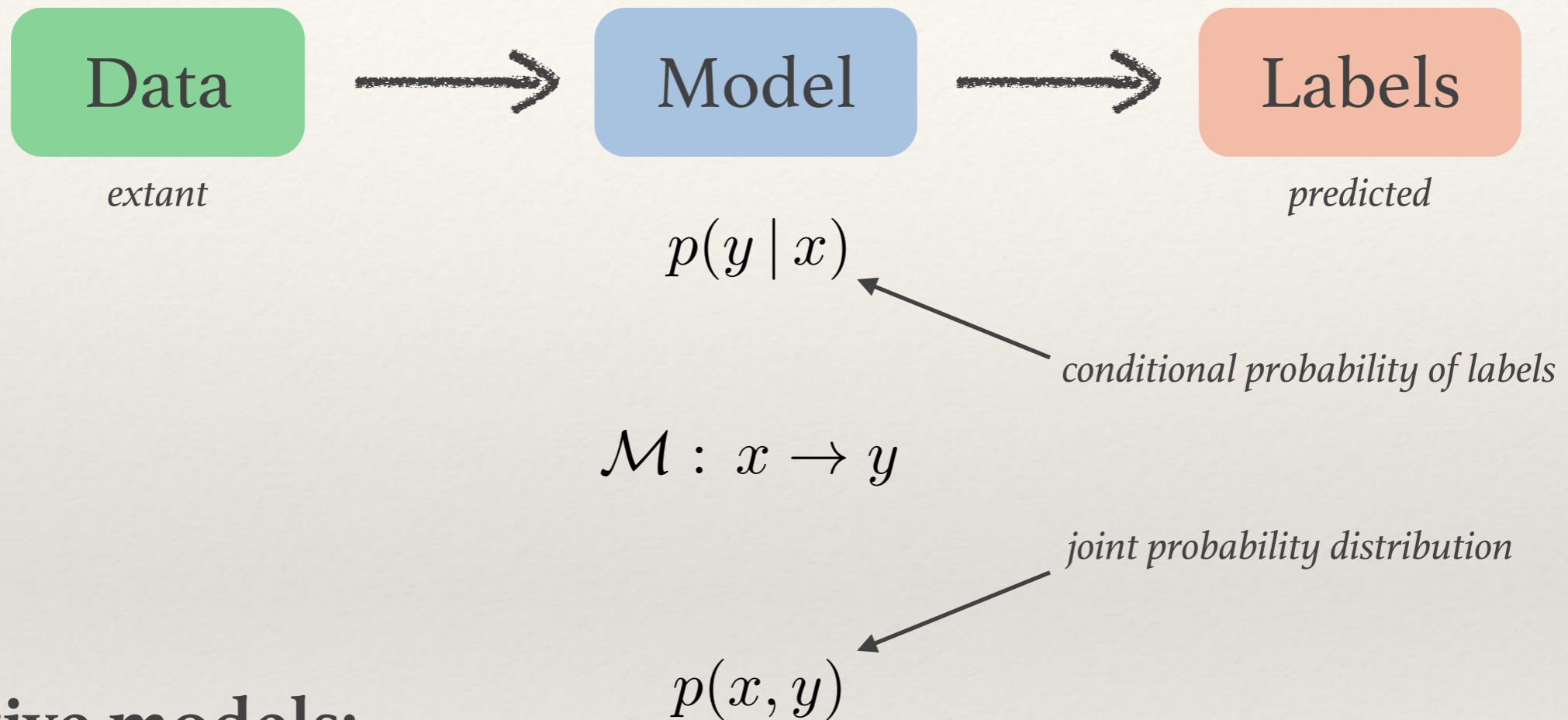
generative models:



↑
learned
↓

How do we do machine learning?

discriminative models:



generative models:



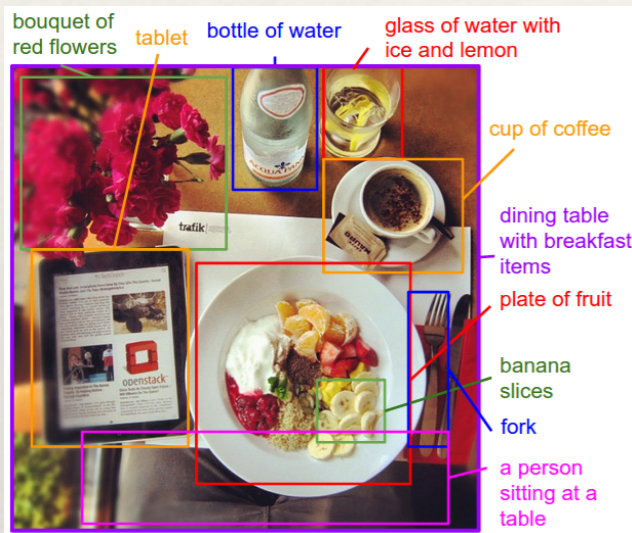
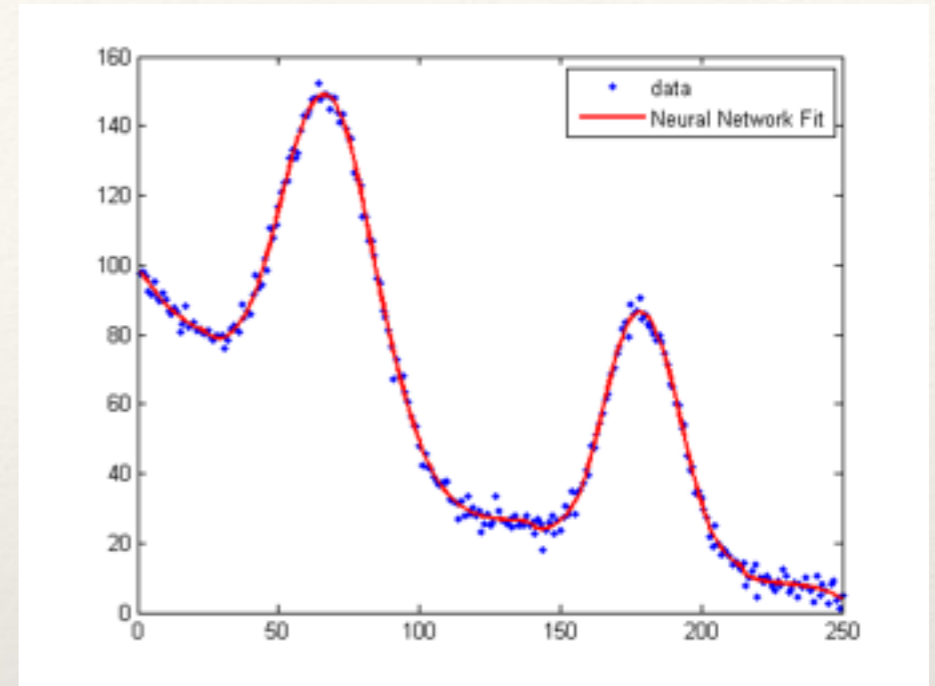
Discriminative Models



categorical discrimination

→ “cat”

regression



semantic labeling

style transfer



AI's will eventually replace us all

This since $\mathcal{F} \in \mathcal{F}$ and $x \in \mathcal{G}$ the diagram

$$\begin{array}{ccc}
 S & \longrightarrow & \\
 \downarrow & & \downarrow \\
 \mathcal{C} & \longrightarrow & \mathcal{O}_{X'} \\
 \text{pt} \downarrow & & \downarrow \\
 & & \mathcal{O}' \\
 & & \downarrow \\
 & & \mathcal{O} \\
 & & \downarrow \\
 & & \text{Spec}(K_\eta)
 \end{array}
 \qquad
 \begin{array}{ccc}
 & & X \\
 & & \downarrow \\
 & & d(\mathcal{O}_{X_{\text{ét}}}, \mathcal{G})
 \end{array}$$

is a limit. Then \mathcal{G} is a finite type and assume S is a flat and \mathcal{F} and \mathcal{G} is a finite type \mathcal{F} . This is of finite type diagrams, and

- the composition of \mathcal{G} is a regular sequence,
- $\mathcal{O}_{X'}$ is a sheaf of rings.

□

Proof. We have see that $X = \text{Spec}(R)$ and \mathcal{F} is a finite type representable by algebraic space. The property \mathcal{F} is a finite morphism of algebraic stacks. Then the cohomology of X is an open neighbourhood of U . □

Proof. This is clear that \mathcal{G} is a finite presentation, see Lemmas ??.

A reduced above we conclude that U is an open covering of \mathcal{C} . The functor \mathcal{F} is a “field”

$$\mathcal{O}_{X,x} \longrightarrow \mathcal{F}_x^{-1}(\mathcal{O}_{X_{\text{ét}}}) \longrightarrow \mathcal{O}_{X,x}^{-1}(\mathcal{O}_{X,x}(\mathcal{O}_{X,x}^{\vee}))$$

is an isomorphism of covering of $\mathcal{O}_{X,x}$. If \mathcal{F} is the unique element of \mathcal{F} such that X is an isomorphism.

The property \mathcal{F} is a disjoint union of Proposition ?? and we can filtered set of presentations of a scheme $\mathcal{O}_{X,x}$ -algebra with \mathcal{F} are opens of finite type over S . If \mathcal{F} is a scheme theoretic image points. □

If \mathcal{F} is a finite direct sum $\mathcal{O}_{X,x}$ is a closed immersion, see Lemma ??.

This is a sequence of \mathcal{F} is a similar morphism.

Generative Models

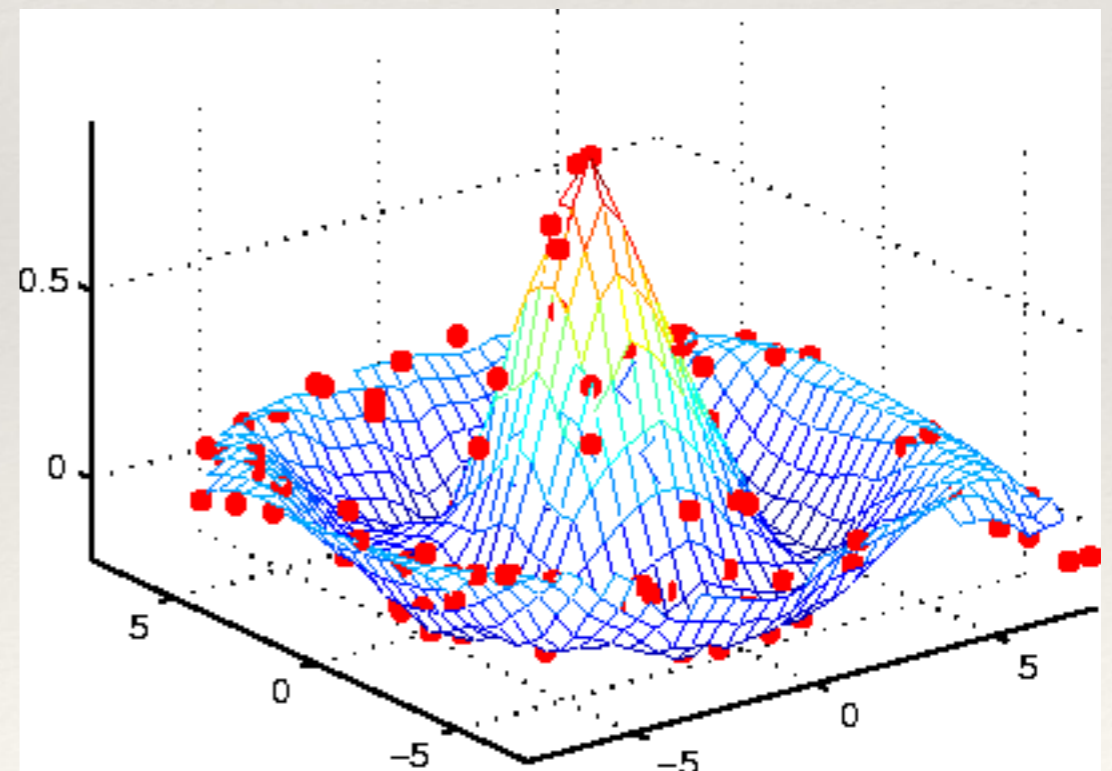
learn a probability distribution: how to sample from data

$$\hat{J}(\Phi)d\Phi \longrightarrow J(\Phi)d\Phi$$

↑ Jacobian (observed via sampling) ↑ learned Jacobian

given points, learn the underlying distribution

physicists do this
all the time



Hopfield Networks

early recurrent neural network (1982)

pairwise connections between nodes

$\{s_i\}$ binary states (-1, +1)

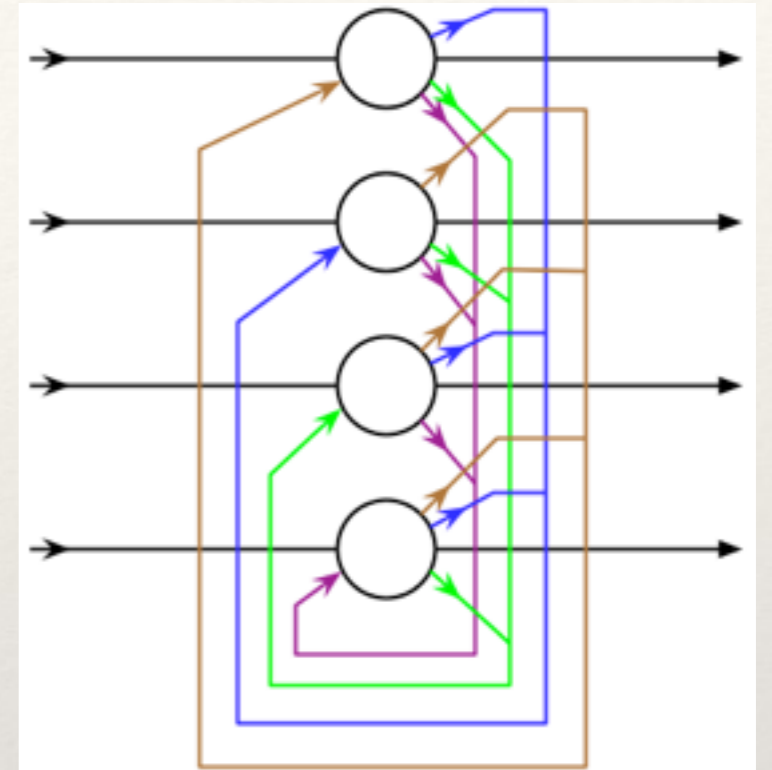
$\{W_{ij}\}$ connection strengths (couplings)

update rule: $s_i = \begin{cases} +1 & \text{if } \sum_j W_{ij} s_j \geq \theta_i \\ -1 & \text{otherwise} \end{cases}$

↑
threshold

energy function: $E = -\frac{1}{2} \sum_{i,j} W_{ij} s_i s_j + \sum_i \theta_i s_i$

under repeated updates, the network converges
to a local minimum in the energy function



The Ising Model

the Hopfield network energy function is similar to an Ising model:

$$E = - \sum_{\langle i j \rangle} J_{ij} s_i s_j - \mu \sum_i h_i s_i \quad \text{Hopfield: all sites 'adjacent'}$$

note that the probability of a given state is dependent on the partition function:

$$P_\beta(s) = \frac{e^{-\beta E(\sigma)}}{Z_\beta}$$

$$Z_\beta = \sum_{\sigma} e^{-\beta E(\sigma)}$$

why Ising models?

simple models that embody the Hebbian learning rule:
neurons that fire together, wire together

concept can be used to store “memories” in the network:
attractor states that are local minima in the energy function

Stochastic Networks

while Hopfield networks are deterministic,
Ising models are probabilistic - like generative models

$$\Delta E_i = E(i \text{ on}) - E(i \text{ off}) \quad \text{energy difference between on/off states}$$

state is activated with probability given by the Boltzmann distribution:

$$\beta \Delta E_i = \ln p_{i \text{ on}} - \ln p_{i \text{ off}}$$

$$p_{i \text{ on}} = \frac{1}{1 + \exp(-\beta \Delta E_i)} = \sigma(-\beta \Delta E_i)$$

this type of network is a **Boltzmann machine**

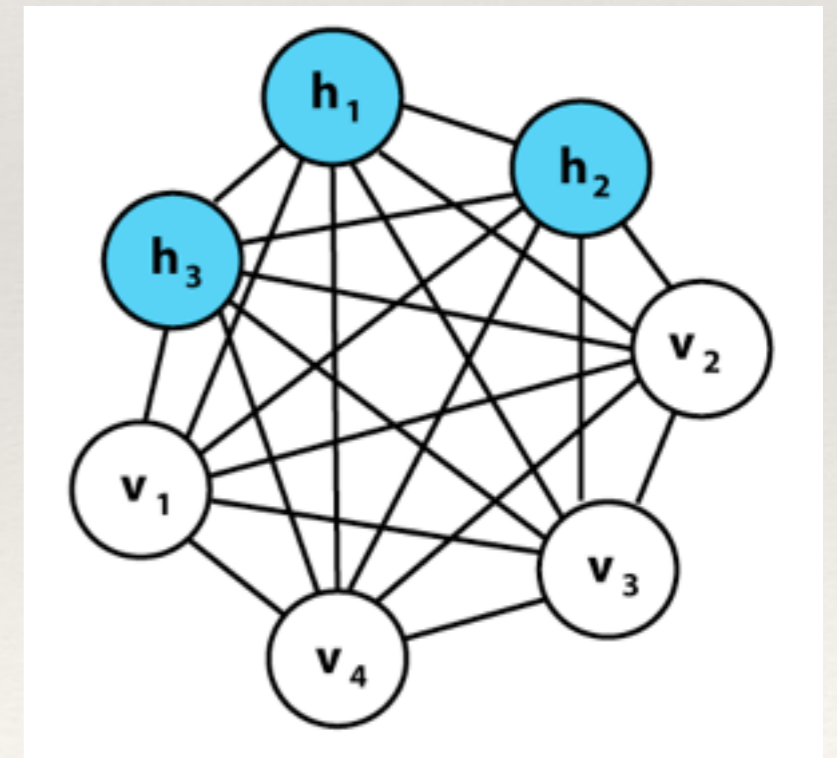
Boltzmann Machines

$$p_{i \text{ on}} = \frac{1}{1 + \exp(-\beta \Delta E_i)} = \sigma(-\beta \Delta E_i)$$

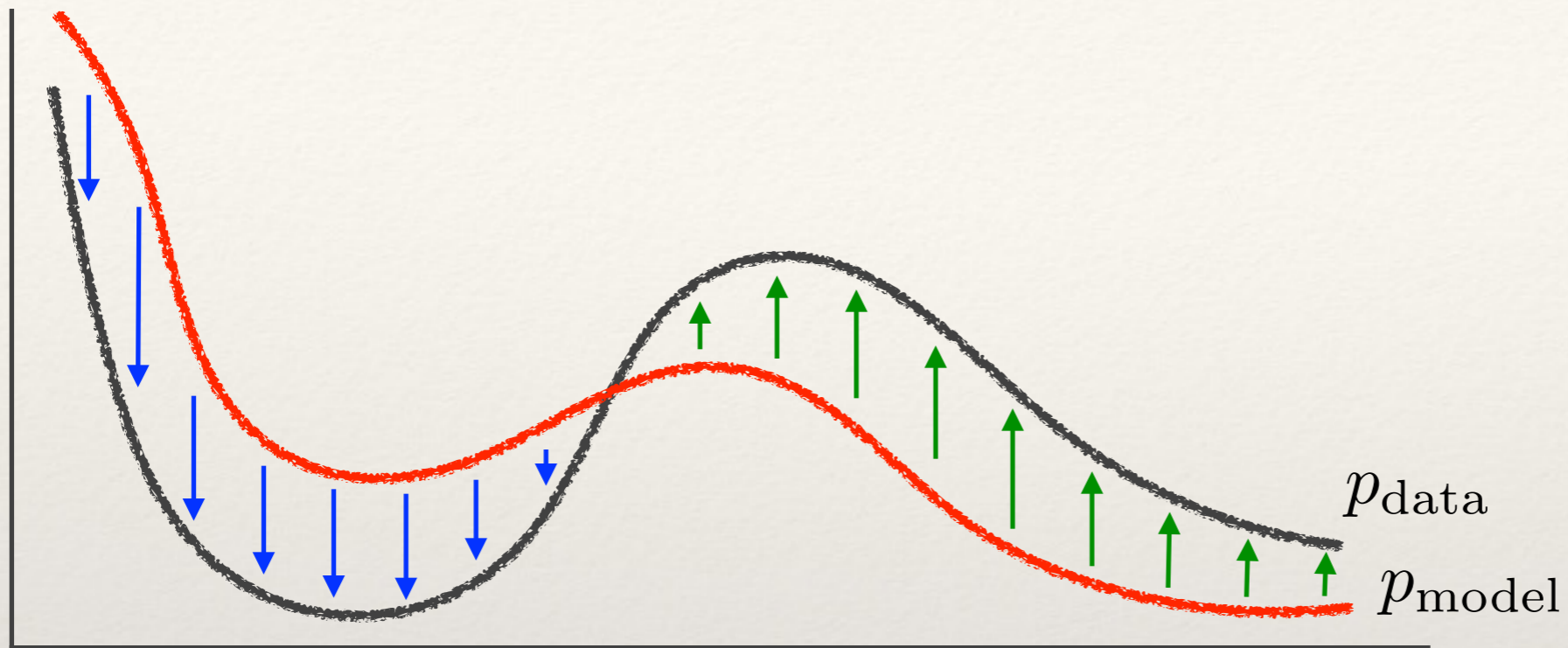
update step: sample sites and set their states according to the Boltzmann distribution; repeat until thermal equilibrium obtained

Structurally, we build a Boltzmann machine from visible (external) and hidden (internal) units

Q: how do we set the weights?
(learning/training)



Training Generative Models



Training: the joint distribution of the model should be adjusted towards the true data distribution

$$\frac{\partial \ln p}{\partial W_{ij}} = (\langle p_{ij} \rangle_{\text{data}} - \langle p_{ij} \rangle_{\text{model}})$$

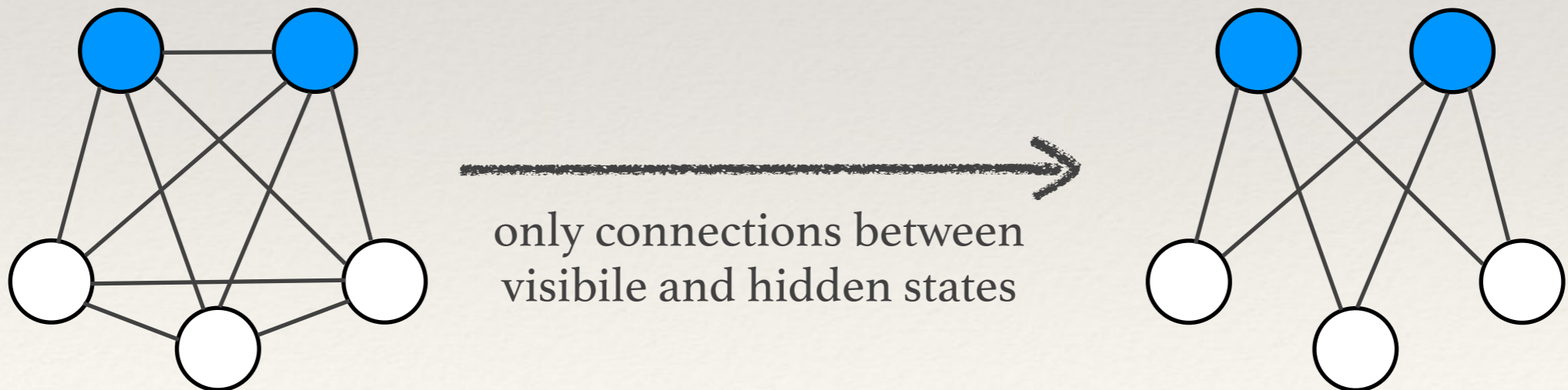
probability of i and j being on

Restricted Boltzmann Machines

training Boltzmann machines is challenging:

- because all states in the network are connected, sampling is extremely time-intensive (single weight updates must propagate through the entire network)
- current training algorithms becomes ineffective beyond small networks

one solution: restricted Boltzmann machines



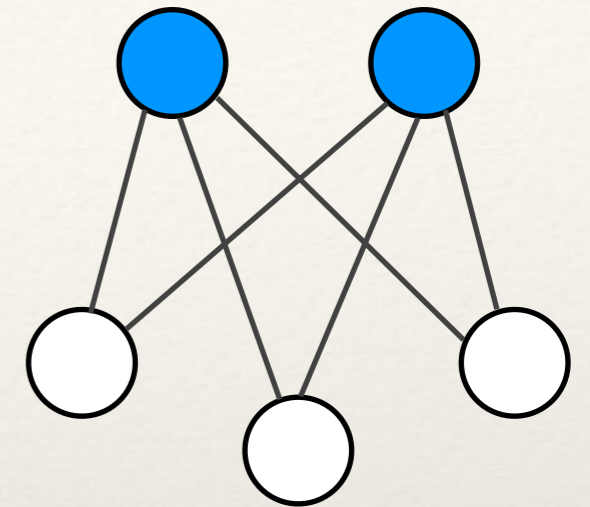
Restricted Boltzmann Machines

RBM energy function:

$$E = - \sum_i a_i v_i - \sum_j b_j h_j - \sum_{i,j} W_{ij} v_i h_j$$

gradient:

$$\frac{\partial \ln p}{\partial W_{ij}} = (\langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{model}})$$

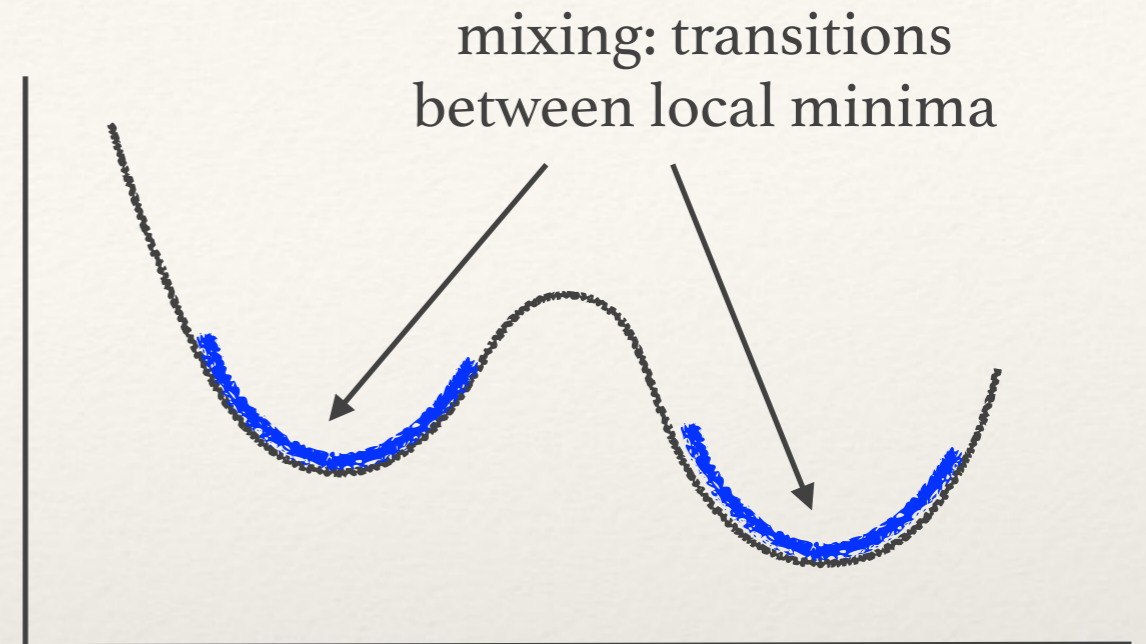


because the visible states are independent and the hidden states are independent, each group can be collectively sampled
- Gibbs sampling

RBM can be effectively trained and used to build expressive models (deep belief networks)

Challenges in Generative Models

- sampling is challenging:
e.g. obtaining *mixing*
- no reliable software packages for training and using RBMs and general Markov process models
- demonstrating a diversity of applications (compete with NNs)



my interests:

- improving sampling methods for generative models using physical systems
- developing useful software tools for experimenting with generative models
- exploring the connection between statistical physics models and ML
- are there good architectures arising from other types of models?

Machine learning is making an enormous impact

The New York Times

Artificial Intelligence Swarms

Silicon Valley on Wings and Wheels

Machine learning is making an enormous impact

The New York Times

*Artificial Intelligence Swarms
Silicon Valley on Wings and Wheels*

And physicists will play a major role in it

WIRED

**MOVE OVER, CODERS—PHYSICISTS WILL
SOON RULE SILICON VALLEY**

Machine learning is making an enormous impact

The New York Times

*Artificial Intelligence Swarms
Silicon Valley on Wings and Wheels*

And physicists will play a major role in it

WIRED

**MOVE OVER, CODERS—PHYSICISTS WILL
SOON RULE SILICON VALLEY**

We even share the same problems

The New York Times

Artificial Intelligence's White Guy Problem

Summary

Machine Learning is just great:

- diverse applications
- exploding interest
- amazing opportunities
- deep roots in statistical physics
- many open questions

Physicists have the tools to make fundamental contributions to machine learning, both inside and beyond physics