

Machine Learning

Lecture 1

School on Data Science in Fundamental Physics, IGFAE/USC, Spain

Dr. Pietro Vischia

pietro.vischia@cern.ch

[@pietrovischia](https://twitter.com/pietrovischia)

<https://vischia.github.io/>



**Supported by project
RYC2021- 033305-I
funded by**

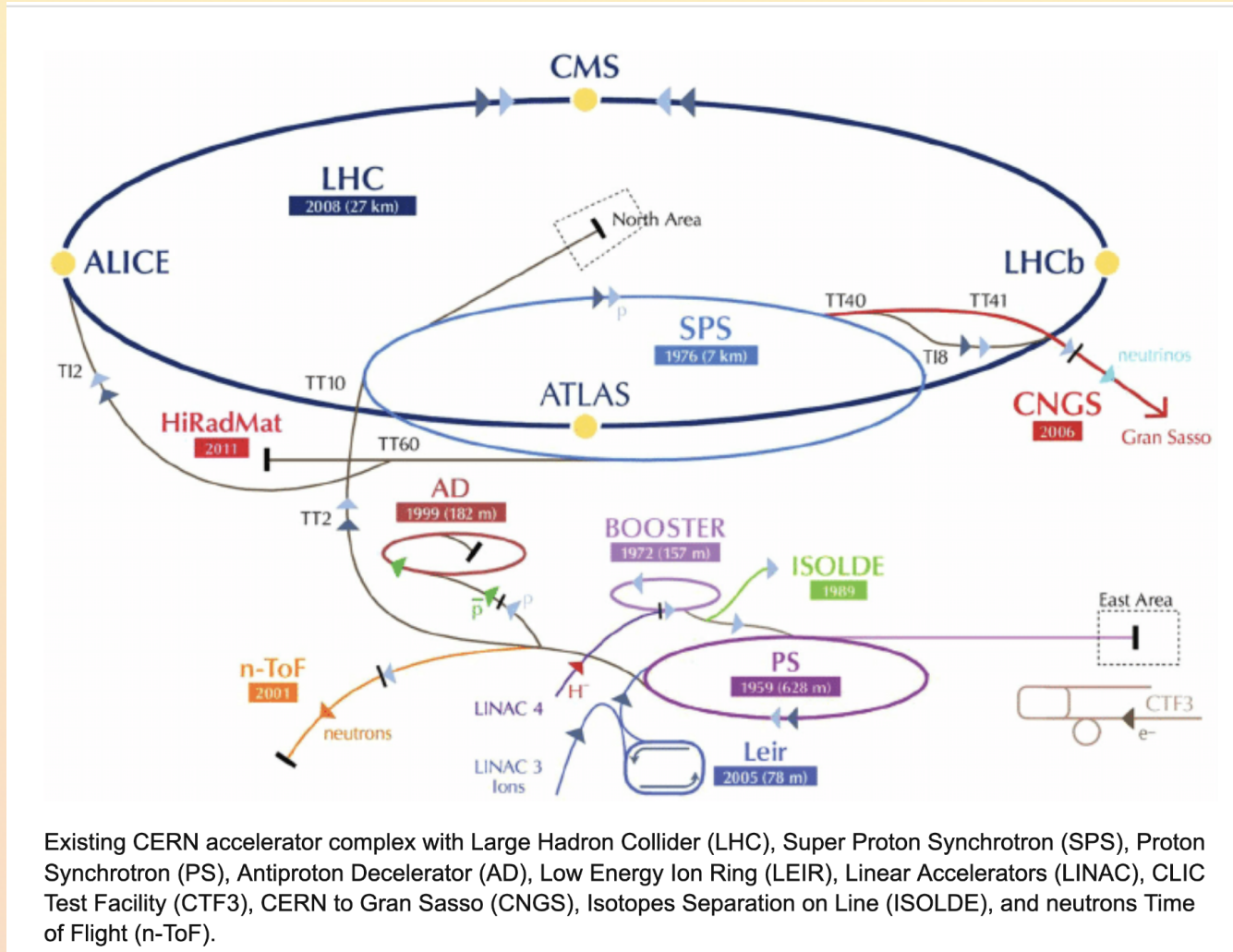


If you are reading this as a web page: have fun! If you are reading this as a PDF:
please visit

https://www.hep.uniovi.es/vischia/persistent/2024-06-03to07_MachineLearningAtDataScienceSchoolIGFAE_vischia_1.html

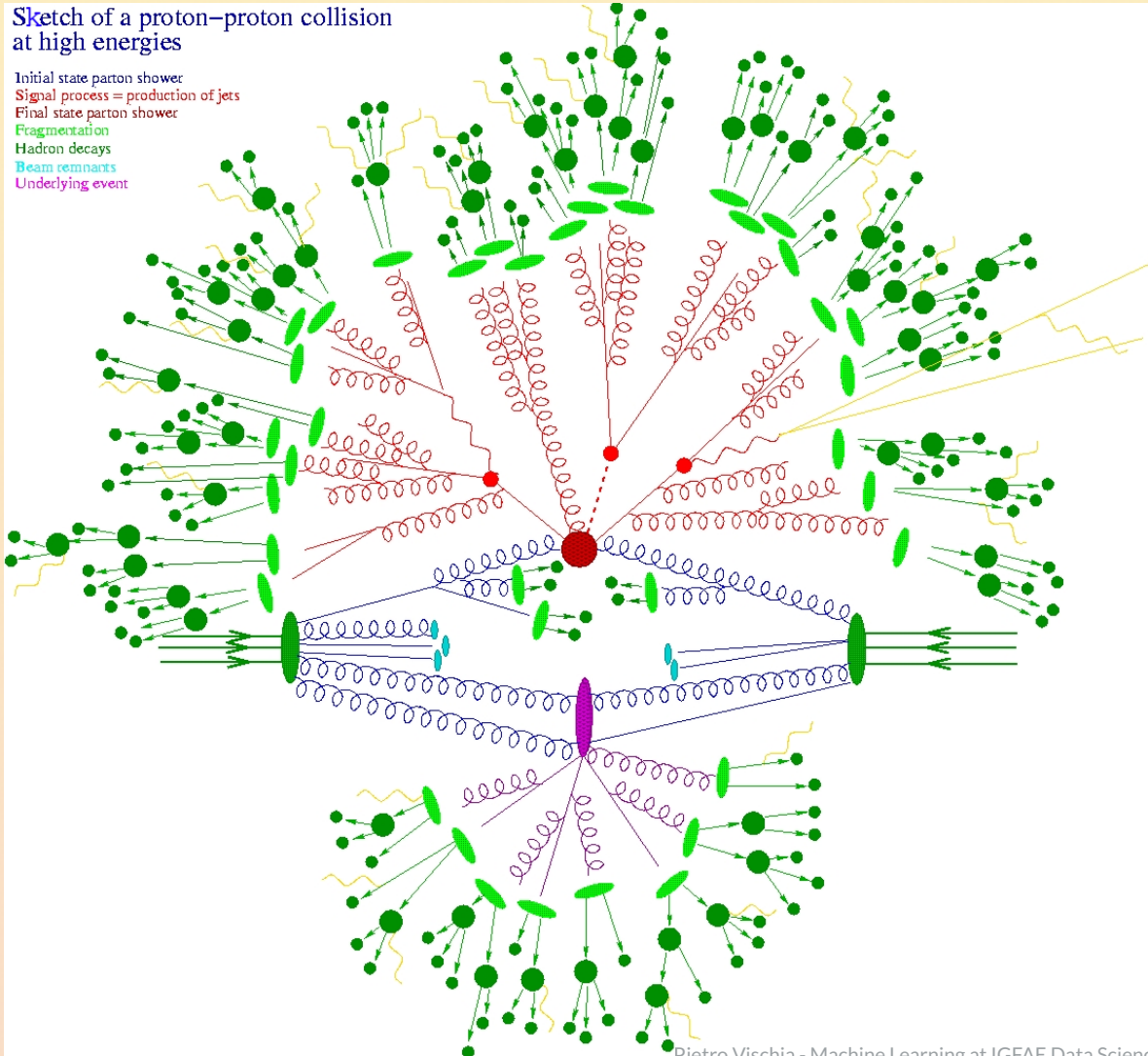
to get the version with working animations

Complex accelerators

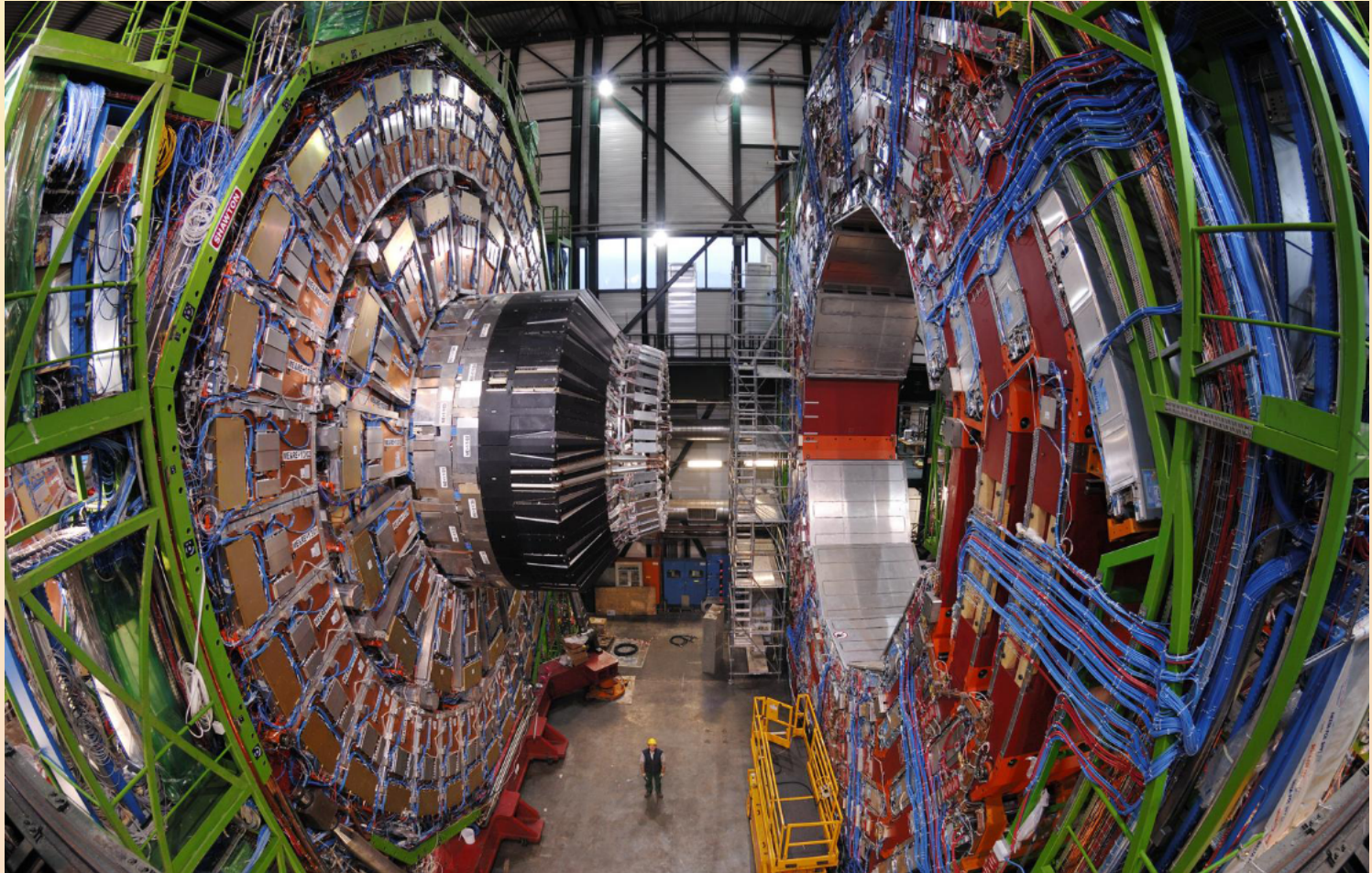


Existing CERN accelerator complex with Large Hadron Collider (LHC), Super Proton Synchrotron (SPS), Proton Synchrotron (PS), Antiproton Decelerator (AD), Low Energy Ion Ring (LEIR), Linear Accelerators (LINAC), CLIC Test Facility (CTF3), CERN to Gran Sasso (CNGS), Isotopes Separation on Line (ISOLDE), and neutrons Time of Flight (n-ToF).

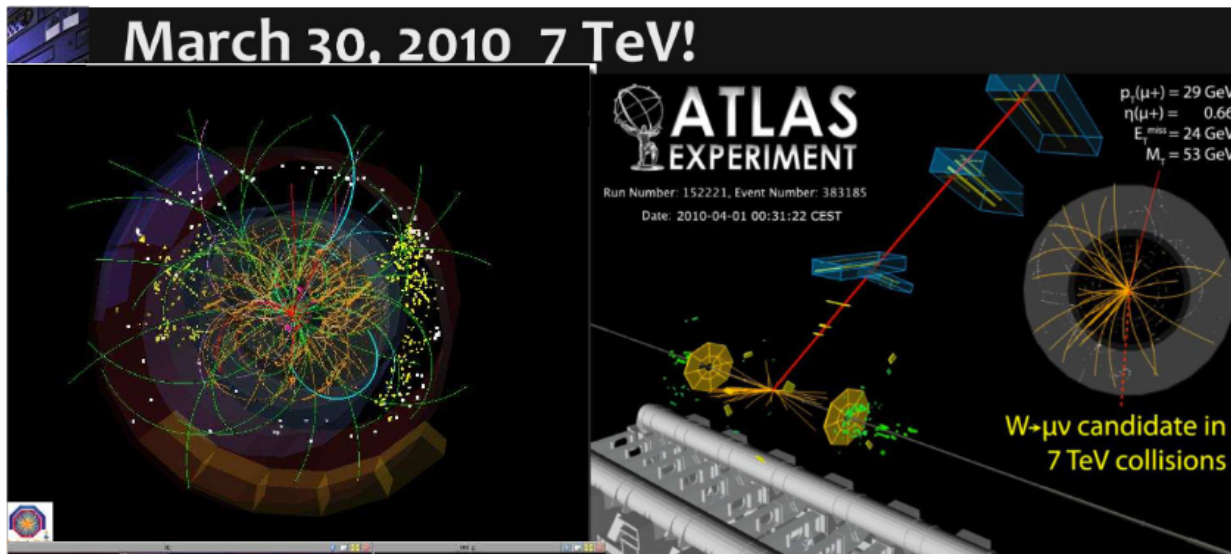
Complex phenomena



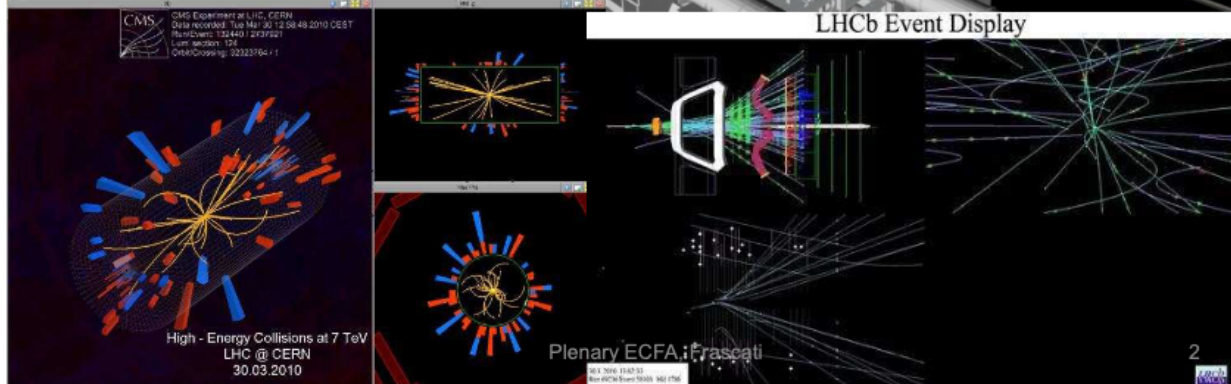
Complex Experiments



What we measure...

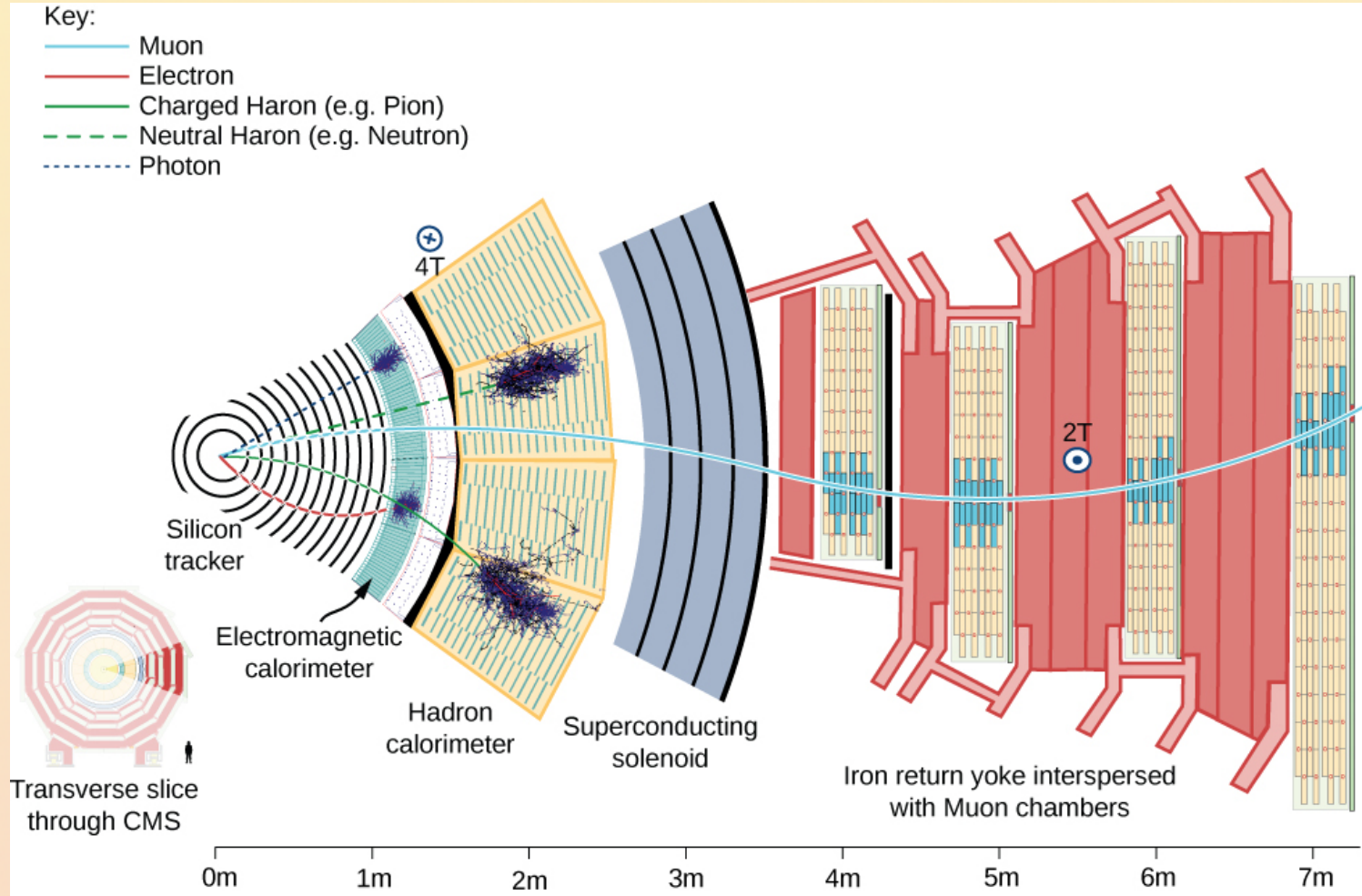


LHCb Event Display

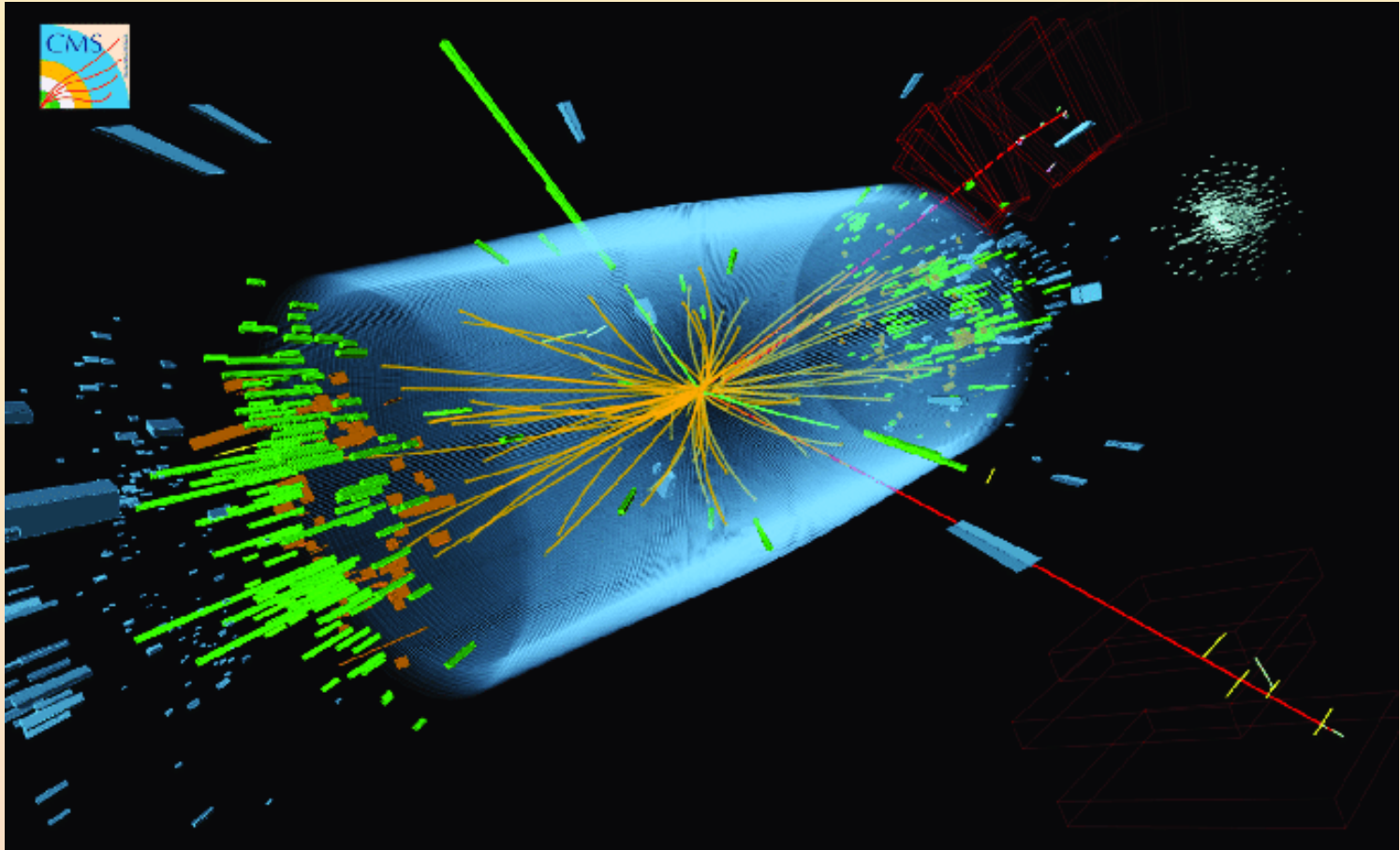


These are the particles that hit the detector

...40 Million Times per Second



Complex Data



Likelihood and information

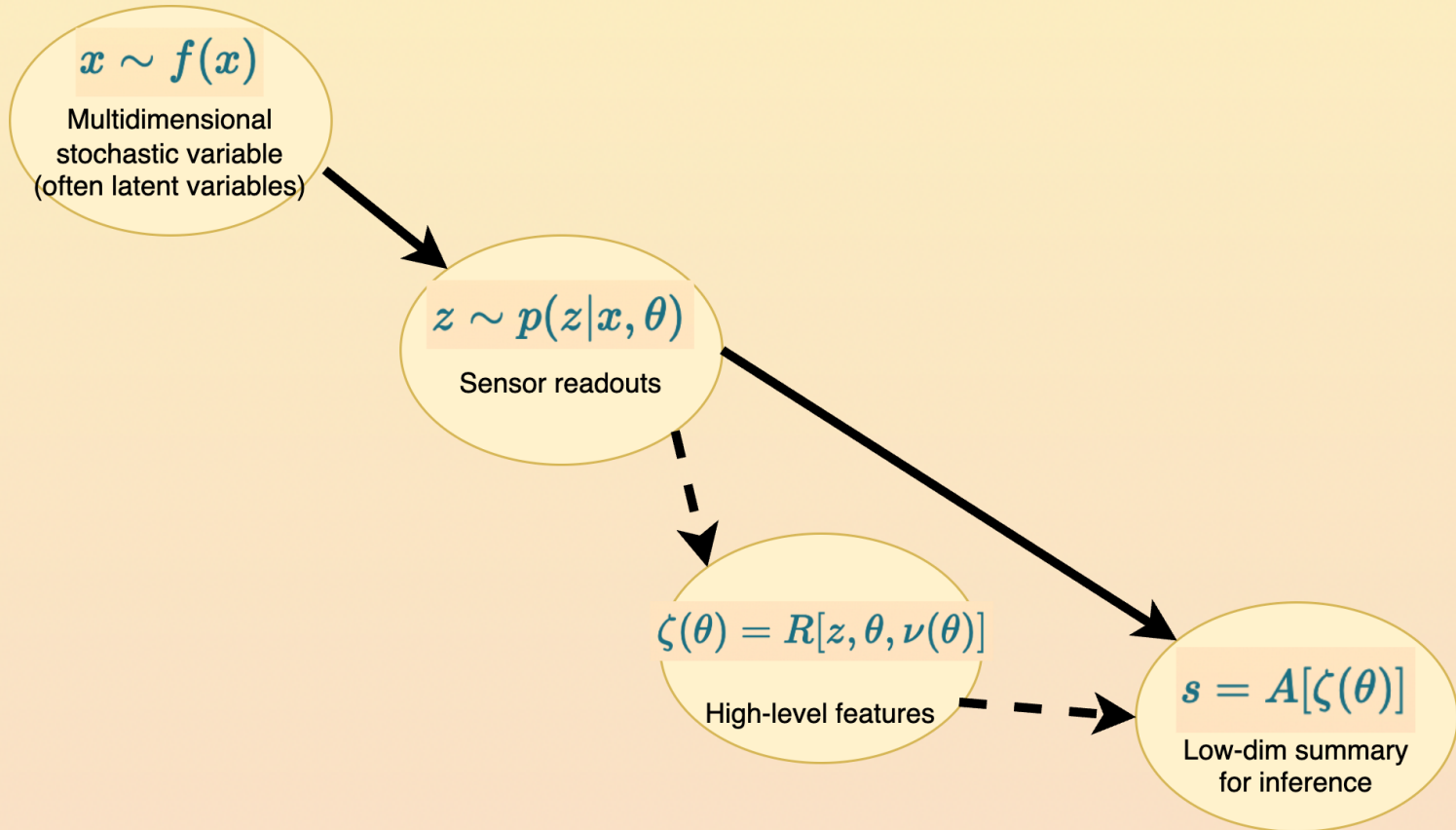
- Data sample X_{obs}

$$\mathcal{L}(X; \theta) := P(X | \theta) |_{X_{obs}}$$

- **The Likelihood Principle:** The likelihood function $L(\vec{x}; \theta)$ contains **all the information** available in the data sample **relevant for the estimation of θ**
 - ✓ Bayesian statistics
 - ✗ Frequentist statistics

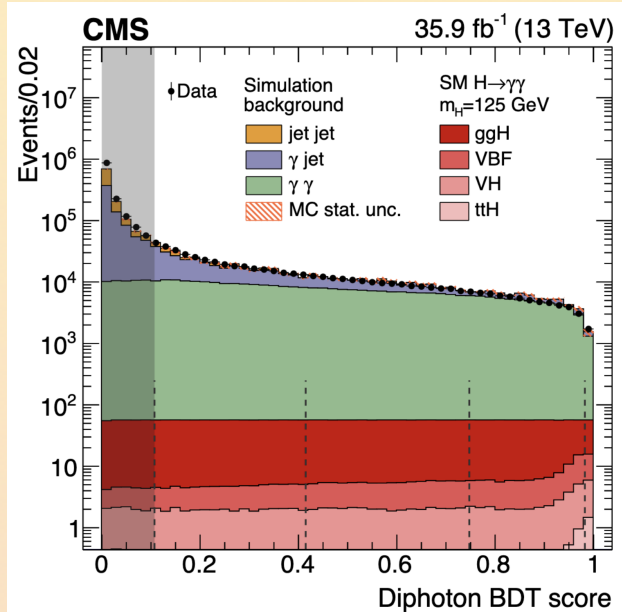
$$I(\theta) = -E \left[\left(\frac{\partial^2}{\partial \theta^2} \ln L(X; \theta) \right)^2 | \theta_{true} \right]$$

Typical analysis pipeline

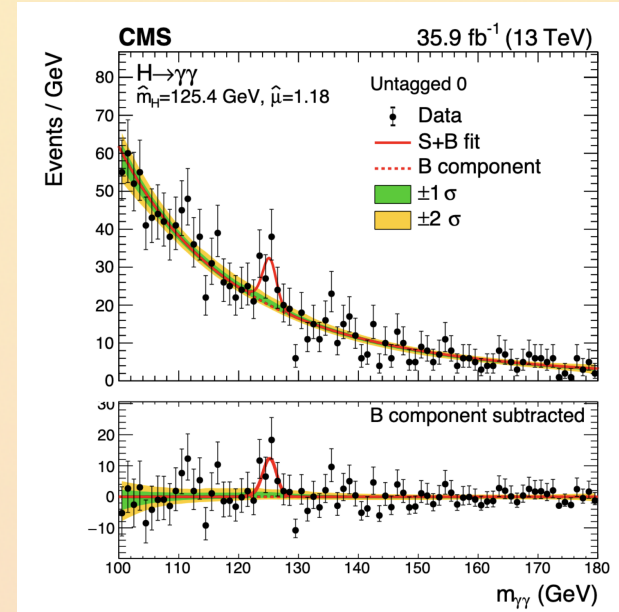


We like low-dim summaries

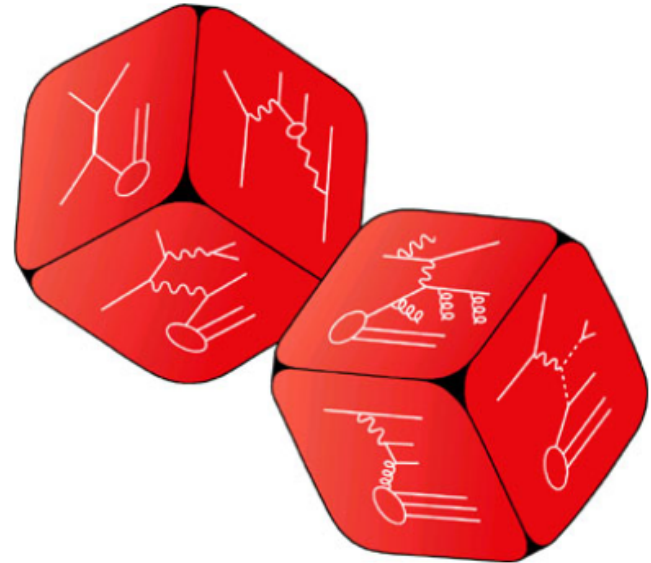
- Discard uninteresting regions



- Physical observable for inference



Quantum mechanics and probability



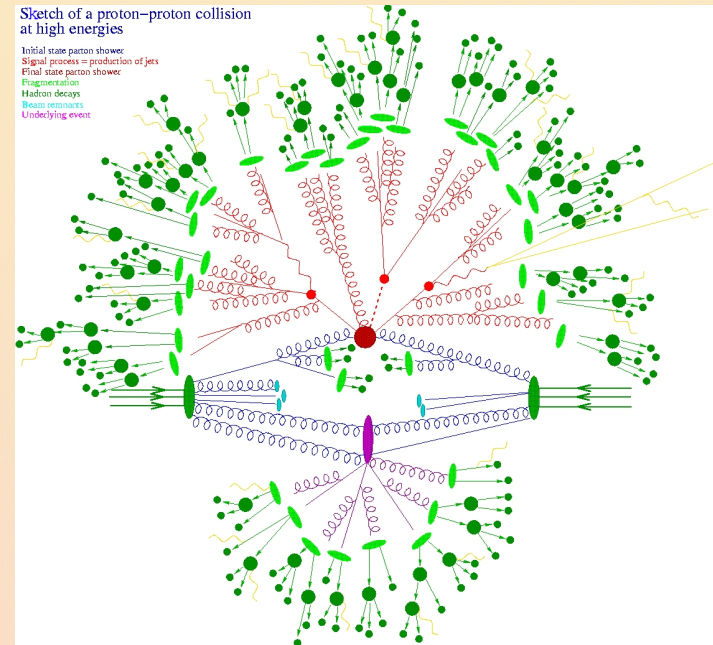
- Quantum mechanical amplitudes as probabilities
- Random numbers as a proxy for quantum mechanical choices

Efficiently sample from models

$$P(\mathbf{x}|\alpha) = \frac{1}{A_\alpha \sigma_\alpha} \int d\Phi(y) \frac{dx_1 dx_2}{x_1 x_2 s} f(x_1) f(x_2) |\mathcal{M}_\alpha(y, x_1, x_2)|^2 W(\mathbf{x}|y) \epsilon_\alpha(y)$$

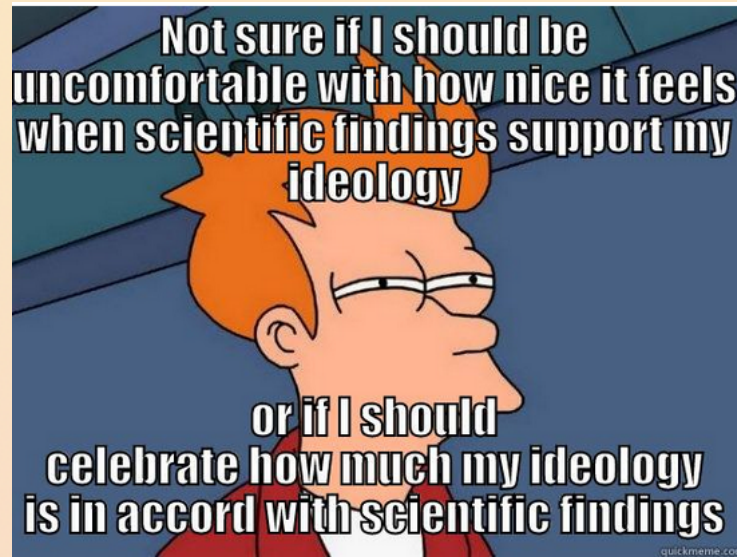
Normalization factor
We collide protons and that is a mess
Linear operator in Hamiltonian formalism, describes the physics process
Detector response, experimental efficiencies

- Difficult integration
- Matrix element computed perturbatively
 - NLO (1st-order) solved in most processes
 - NNLO current challenge
- Costly MonteCarlo simulators
 - Hundreds-dimensional phase space



Generators as vehicles of ideology

- Study complex multiparticle physics (both experimentally and theoretically)
- Flexibility in physical quantities that can be addressed
 - Predict event rates and topologies (estimate feasibility)
 - Simulate possible backgrounds (devise analysis strategies)
 - Study detector requirements (can optimize detector/trigger design)
 - Study detector imperfections (evaluate acceptance corrections)



Monte Carlo simulations

- We want to generate events as they are in Nature
 - average and fluctuations
 - random choices

$$\sigma_{final\ state} = \sigma_{hard\ process} P_{tot, hard\ process \rightarrow final\ state}$$

$$\text{where } P_{tot} = P_{res} P_{ISR} P_{FSR} P_{MPI} P_{remnants} P_{hadronization} P_{decays}$$

$$\text{and } P_i = \prod_j P_{ij} = \prod_i \prod_k P_{ijk} = \dots$$

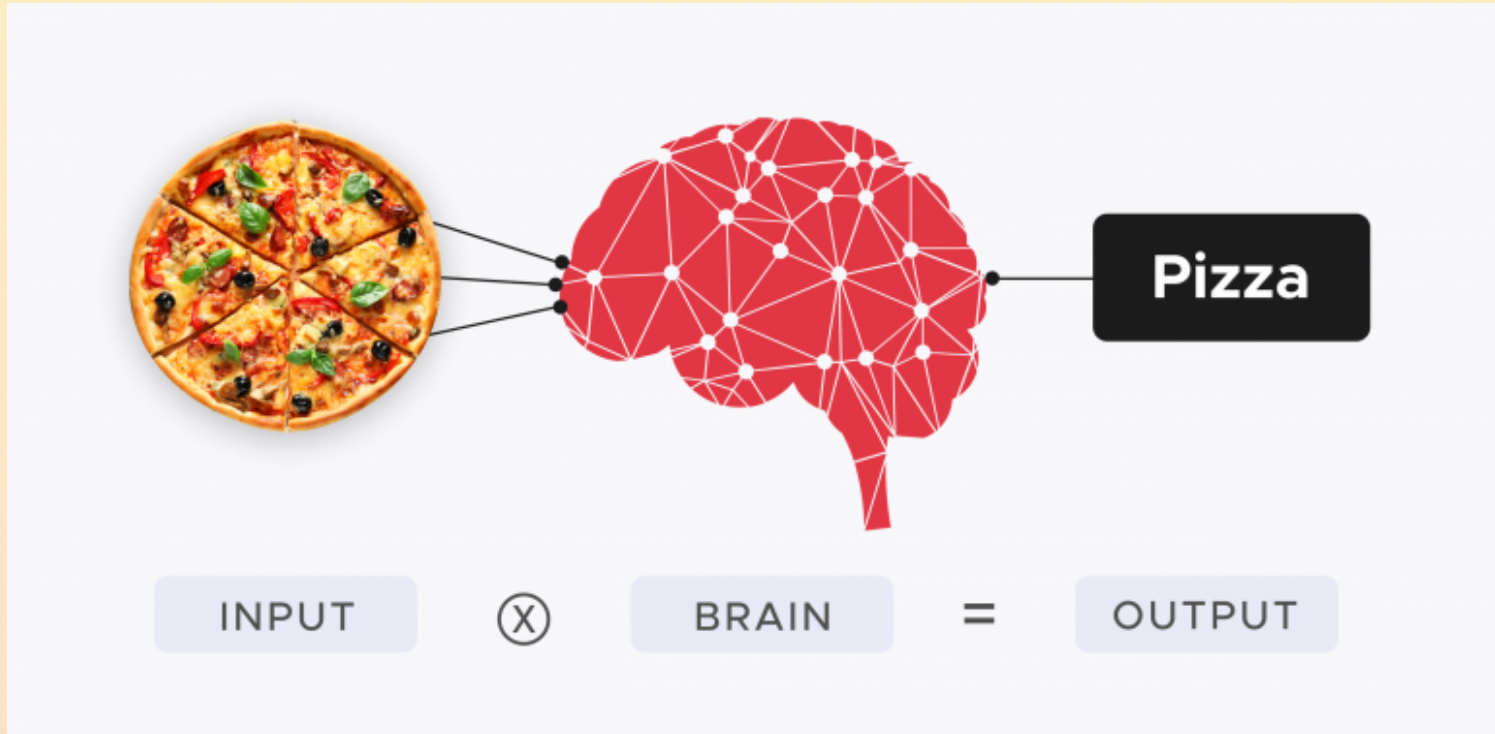
- For each event, $\mathcal{O}(10)$ random choices (flavour, mass, momentum, spin, lifetime, ...)
- Typical LHC event: ~ 100 charged particles and ~ 200 neutral particles
- Several thousand choices

Likelihood intractability in HEP

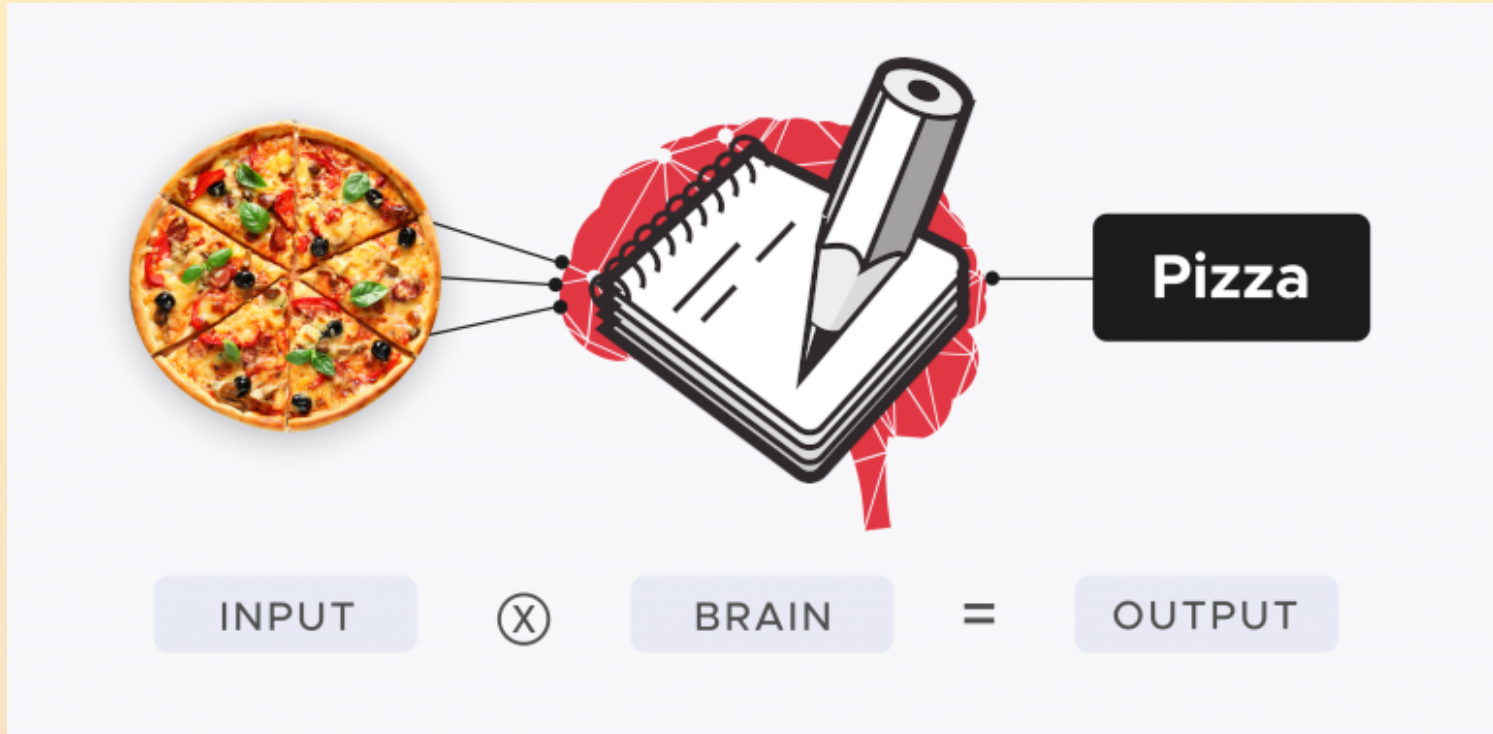
$$p(x|\theta) = \int dz p(x, z|\theta) = \int dz p_x(x|\theta, z) \prod_i p_i(z_i|\theta, z_{<i})$$

- No access to the likelihood, or latent variables
 - Matrix element
 - Parton shower
 - Detector simulation
- Can always generate Monte Carlo samples from $x \sim p(x|\theta)$
- Histograms are likelihood-free!!!
 - Count events, assume Poisson per bin, global likelihood as product

Brain activity...



...approximated...



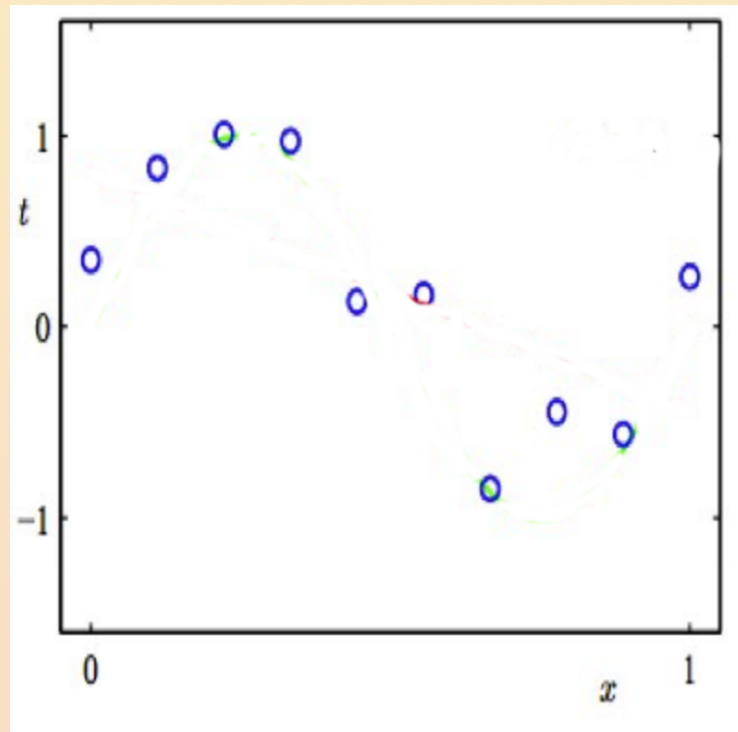
...using computers



Understanding Data

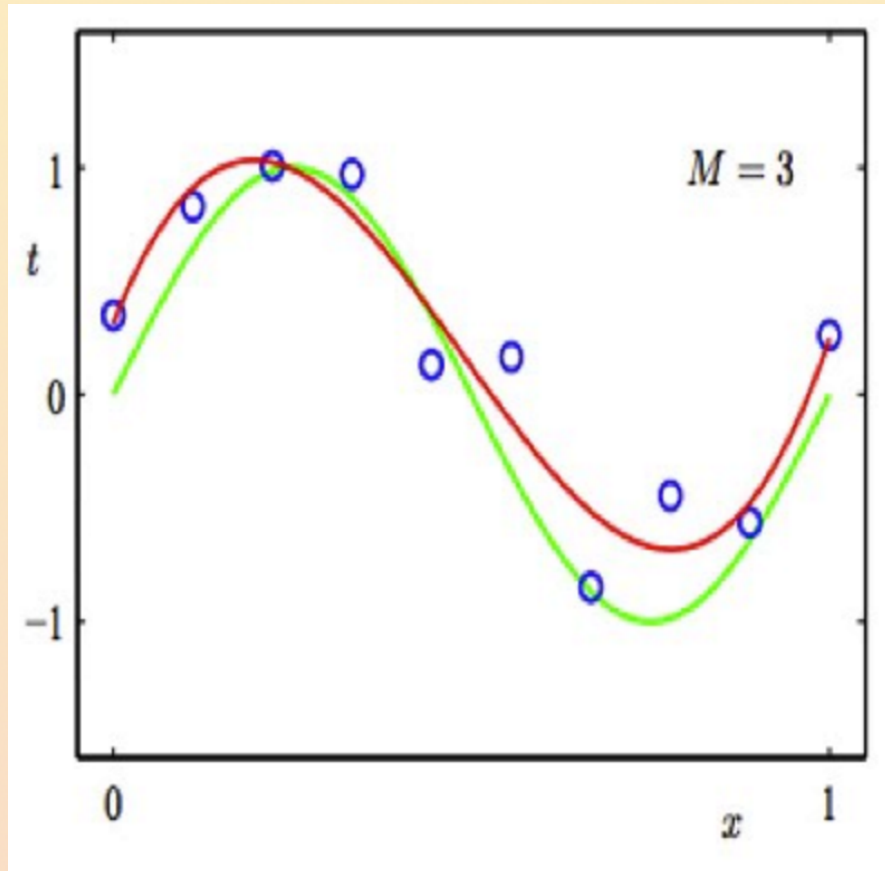
Vast amounts of data are being generated in many fields, and the statistician's job is to make sense of it all: to extract important patterns and trends, and understand "what the data says." We call this learning from data.

(Hastie, Tibshirani, Friedman, Springer 2017)



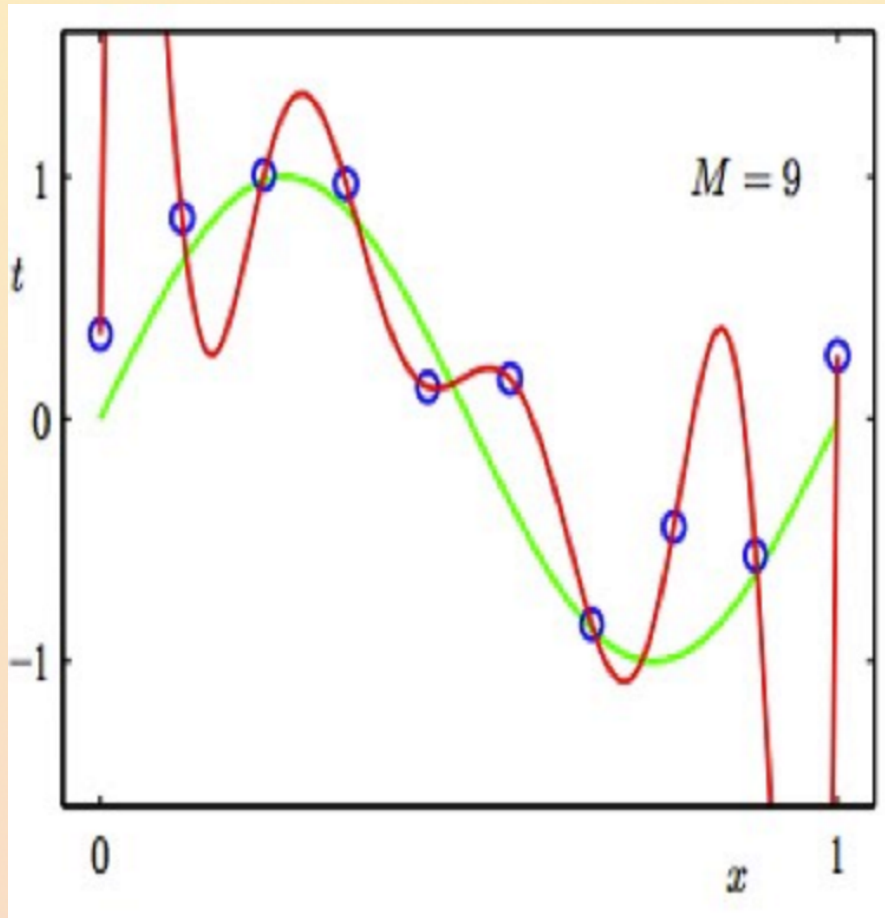
Functions Describe the World

- Interpolation

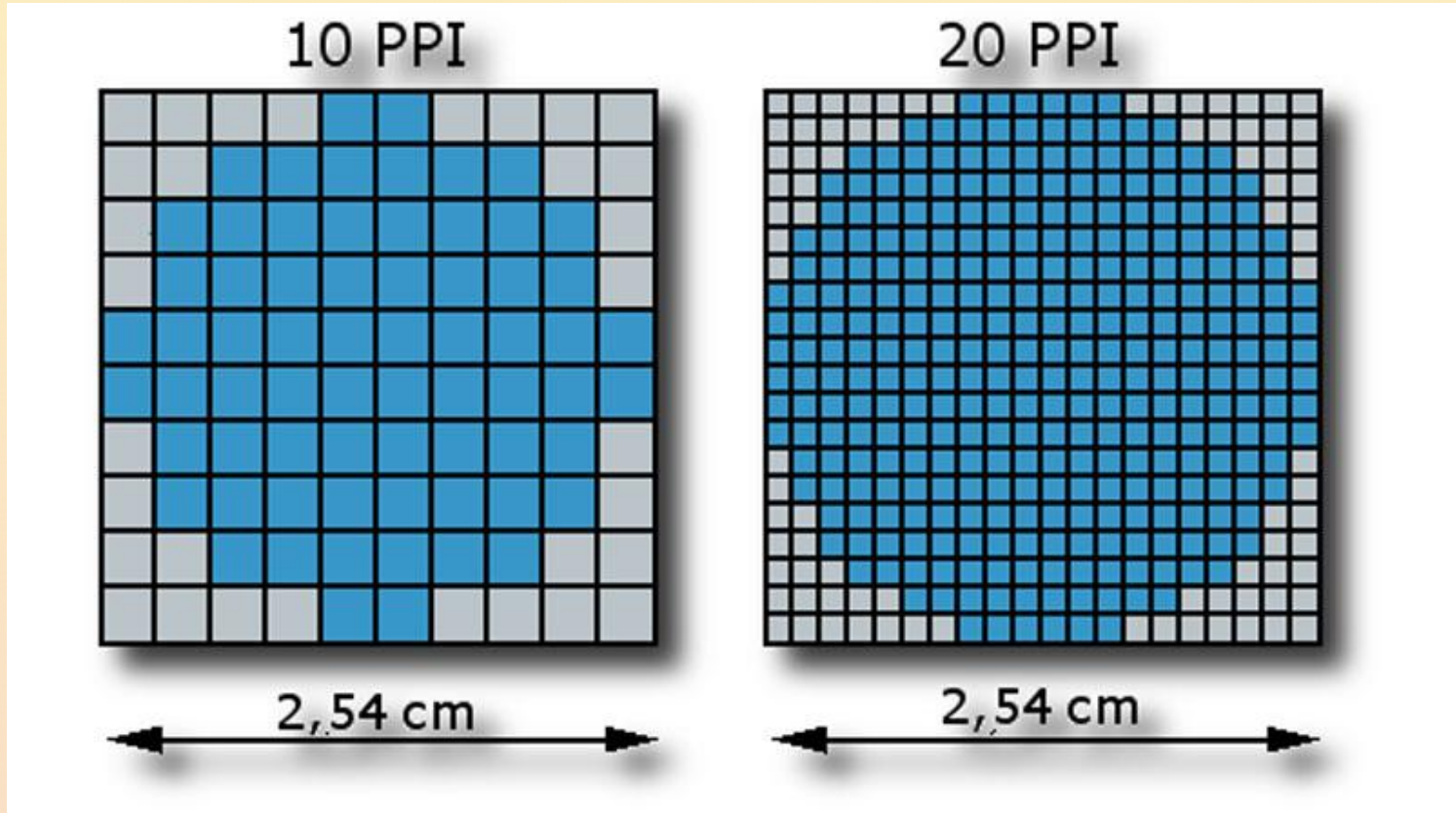


Sometimes too well

- Generalization



Think in Millions of Dimensions



Easy or difficult?

```
[vischia@lxplus797 ~]$ python dump_cpTree.py
Info in <TCanvas::MakeDefCanvas>: created default TCanvas with name c1
*****
* Row * Instance * Lep1_pt.L * Lep2_pt.L * Lep3_pt.L * Lep1_eta. * Lep2_eta. * Lep3_eta. * Lep1_phi. * Lep2_phi. * Lep3_phi. * met.met * met_phi.m * weight_CP * HTT_score *
*****
* 56350 * 0 * 280.82870 * 93.600341 * 59.040779 * -1.360351 * 0.7332763 * -0.819335 * -2.826171 * 0.7969970 * -3.100097 * 118.03358 * 2.5146484 * 40.820312 * 0.8255668 *
* 56350 * 1 * 280.82870 * 93.600341 * 59.040779 * -1.360351 * 0.7332763 * -0.819335 * -2.826171 * 0.7969970 * -3.100097 * 118.03358 * 2.5146484 * 40.820312 * 0.8255668 *
* 56350 * 2 * 280.82870 * 93.600341 * 59.040779 * -1.360351 * 0.7332763 * -0.819335 * -2.826171 * 0.7969970 * -3.100097 * 118.03358 * 2.5146484 * 40.820312 * 0.8255668 *
* 56350 * 3 * 280.82870 * 93.600341 * 59.040779 * -1.360351 * 0.7332763 * -0.819335 * -2.826171 * 0.7969970 * -3.100097 * 118.03358 * 2.5146484 * 40.820312 * 0.8255668 *
* 79791 * 0 * 67.791183 * 42.294036 * 17.310911 * -1.846923 * 1.4458007 * -0.762207 * -0.628540 * -2.910644 * 2.9912109 * 8.8895807 * -1.773925 * 94.398437 * -99 *
*****
=> 5 selected entries
-----
* Row * Instance * Lep1_pt.L * Lep2_pt.L * Lep3_pt.L * Lep1_eta. * Lep2_eta. * Lep3_eta. * Lep1_phi. * Lep2_phi. * Lep3_phi. * met.met * met_phi.m * weight_CP * HTT_score *
*****
* 58751 * 0 * 86.053443 * 32.593345 * 13.077508 * -1.346435 * 2.0512695 * -0.503906 * -0.392944 * -2.925781 * 0.9309082 * 17.544691 * 0.1735229 * 3.478e-08 * 0.9726203 *
* 58751 * 1 * 86.053443 * 32.593345 * 13.077508 * -1.346435 * 2.0512695 * -0.503906 * -0.392944 * -2.925781 * 0.9309082 * 17.544691 * 0.1735229 * 3.478e-08 * 0.9726203 *
* 58751 * 2 * 86.053443 * 32.593345 * 13.077508 * -1.346435 * 2.0512695 * -0.503906 * -0.392944 * -2.925781 * 0.9309082 * 17.544691 * 0.1735229 * 3.478e-08 * 0.9726203 *
* 58751 * 3 * 86.053443 * 32.593345 * 13.077508 * -1.346435 * 2.0512695 * -0.503906 * -0.392944 * -2.925781 * 0.9309082 * 17.544691 * 0.1735229 * 3.478e-08 * 0.9726203 *
* 58751 * 4 * 86.053443 * 32.593345 * 13.077508 * -1.346435 * 2.0512695 * -0.503906 * -0.392944 * -2.925781 * 0.9309082 * 17.544691 * 0.1735229 * 3.478e-08 * 0.9726203 *
* 58751 * 5 * 86.053443 * 32.593345 * 13.077508 * -1.346435 * 2.0512695 * -0.503906 * -0.392944 * -2.925781 * 0.9309082 * 17.544691 * 0.1735229 * 3.478e-08 * 0.9726203 *
*****
=> 6 selected entries
```

Signal

Background



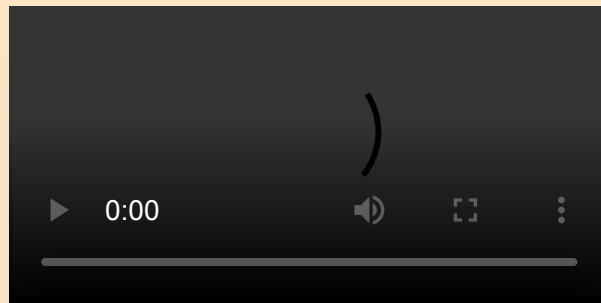
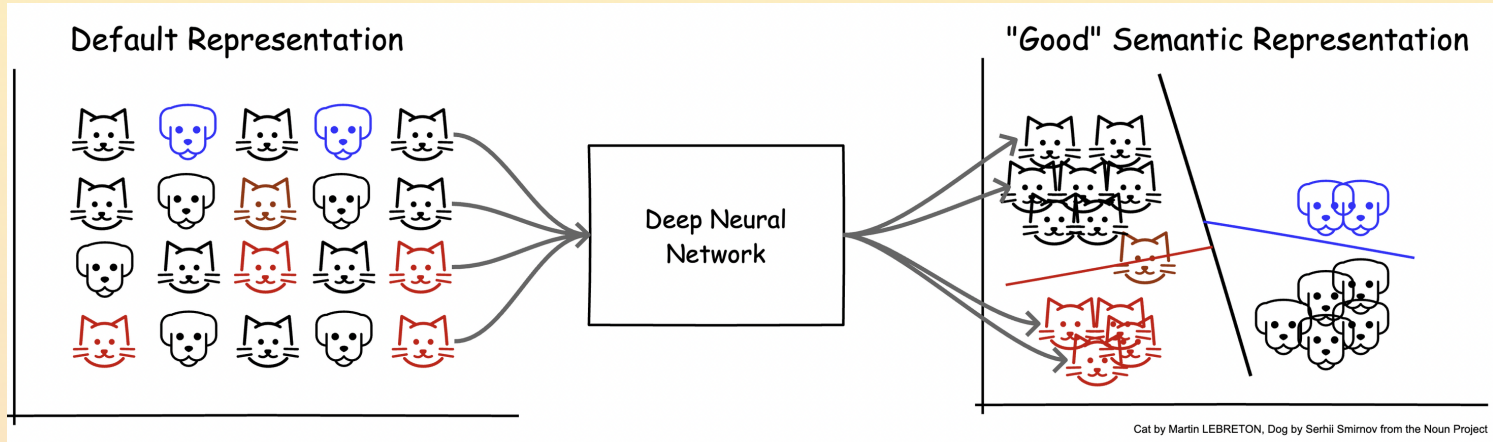
Easy or difficult?



Mapping Improves Understanding

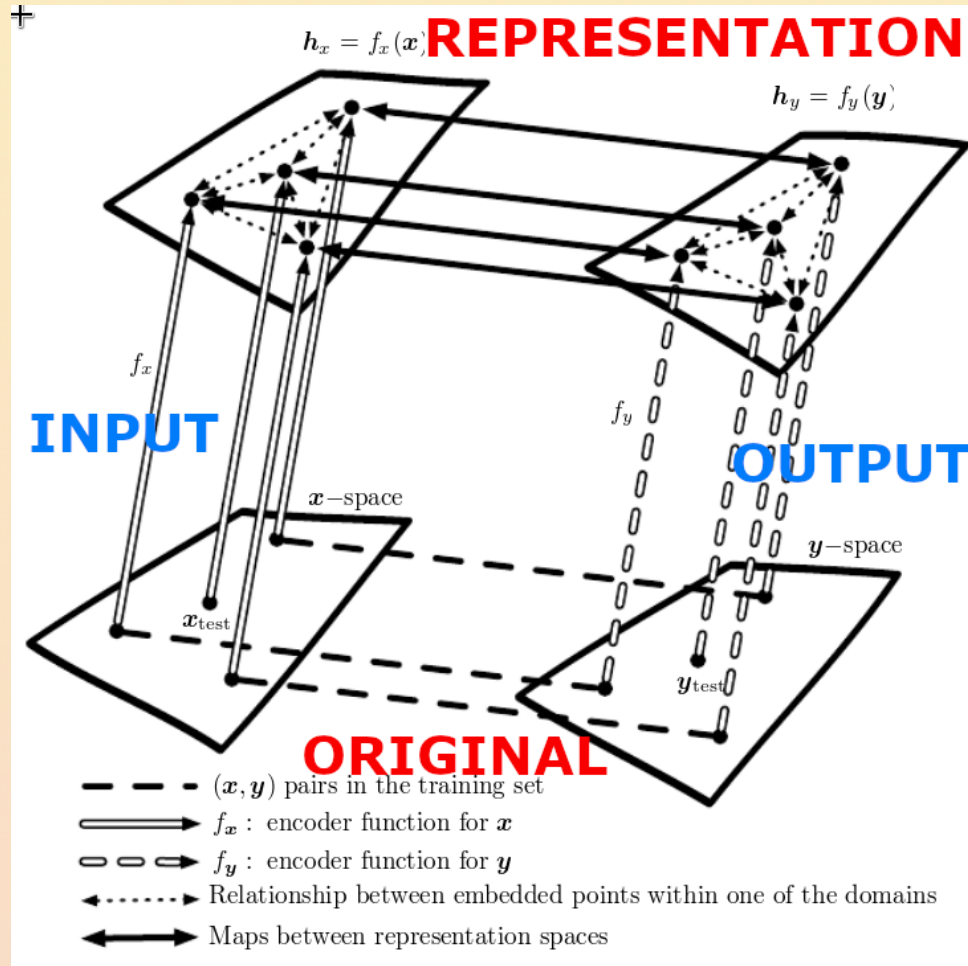


Representations Make Tasks Easier



Learn Representations

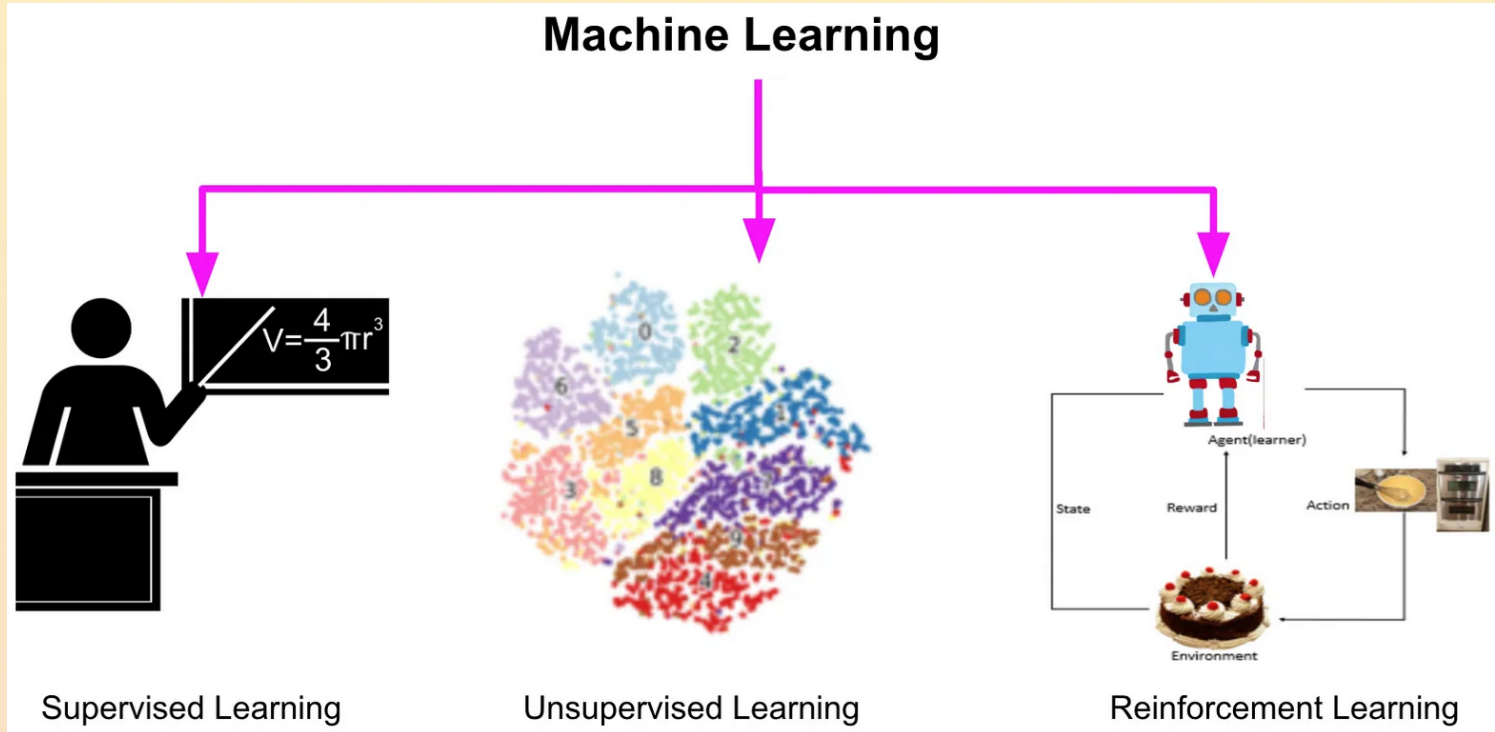
- Like AdS/CFT, but actually demonstrated 😊



Learning from data:

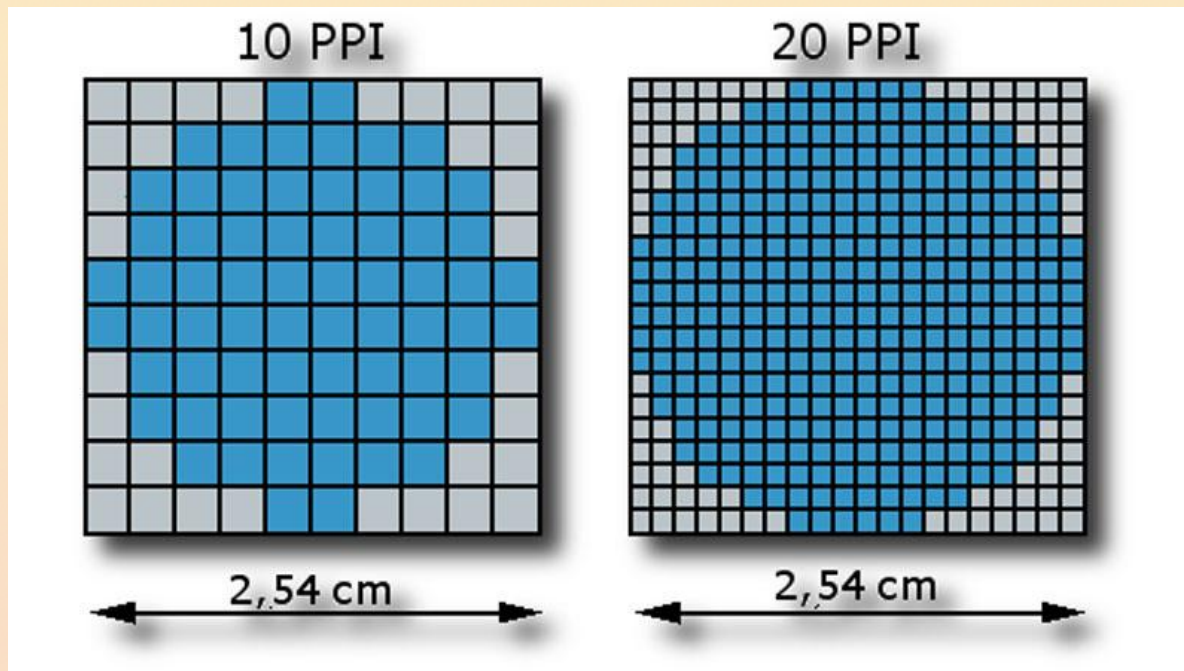
the mathematical formalism

Learn in different ways



Input space

- \mathcal{X} : a high-dimensional input space
 - The challenges come from the high dimensionality!
- If all dimensions are real-valued, \mathbb{R}^d
 - For square images of side \sqrt{d} , space is $\mathcal{X} = \mathbb{R}^d$, and $d \sim \mathcal{O}(10^6)$



Data probability distribution

- ν : unknown data probability distribution
 - We can sample from it to obtain an arbitrary amount of data points
 - We are not allowed to use any analytic information about it in our computations

The target function

- $f^* : \mathcal{X} \rightarrow \mathbb{R}$, unknown target function
 - In case of multidimensional output to a vector of dimension k , $f^* : \mathcal{X} \rightarrow \mathbb{R}^k$
 - Some loose assumptions (e.g. square-integrable with respect to the ν measure, i.e. finite moments, bounded...)

The loss functional

- $L[f] = \mathbb{E}_{\nu} \left[l\left(f(x), f^*(x)\right) \right]$
 - The metric that tells us how good our predictions are
- The function $l(\cdot, \cdot)$ is a given expression, e.g. regression loss, logistic loss, etc
 - In this lecture, typically it is the L^2 norm: $\|f - f^*\|_{L^2(\mathcal{X}, \nu)}$

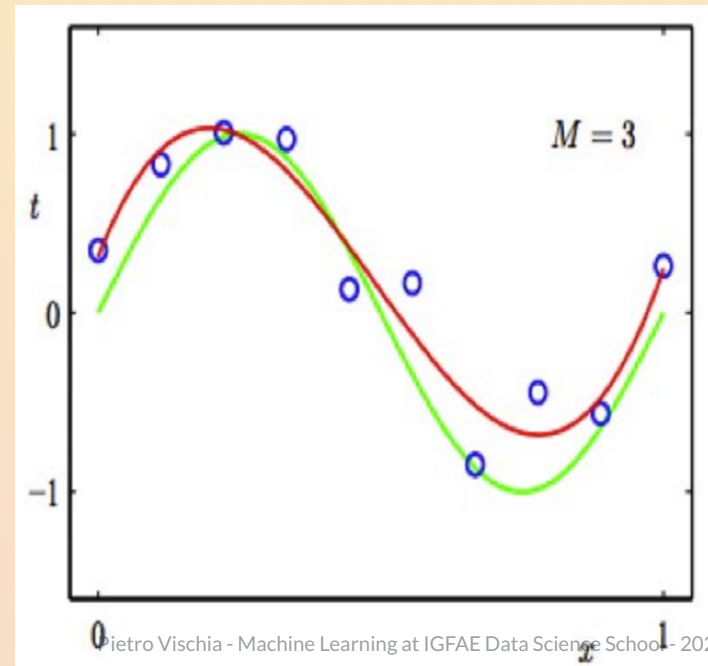
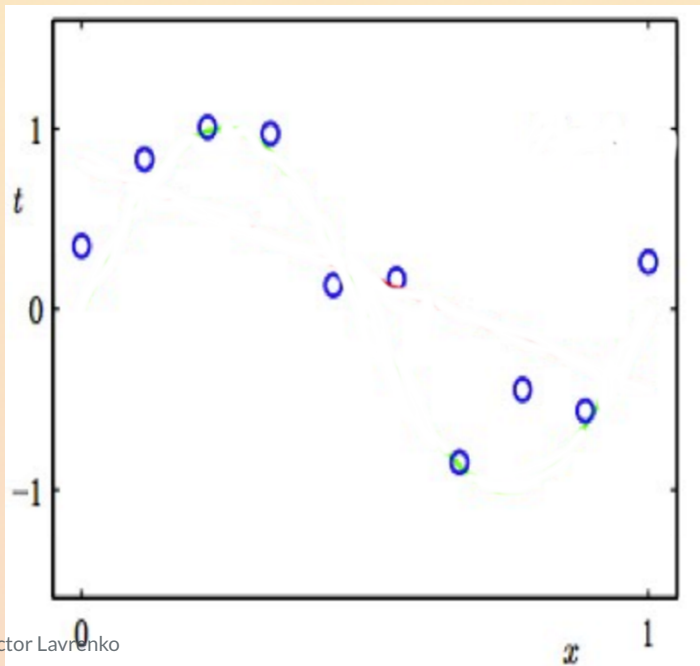
Loss function comes from inference

- Decision-theoretic approach (C.P. Robert, "The Bayesian Choice")
 - \mathcal{X} : observation space
 - Θ : parameter space
 - \mathcal{D} : decision (action) space
- **Statistical inference** take a decision $d \in \mathcal{D}$ related to parameter $\theta \in \Theta$ based on observation $x \in \mathcal{X}$, under $f(x|\theta)$
 - Typically, d consists in estimating $h(\theta)$ accurately

$$U(\theta, d) = \mathbb{E}_{\theta, d} \left[U(r) \right]$$

Learning goal

- Goal: predict f^* from a finite i.i.d. sample of points sampled from ν
- Sample $\{x_i, f^*(x_i)\}_{i=1, \dots, n}, x_i \sim \nu$
 - For each of the points x_i , we know the value of the unknown function (our true labels)
 - We want to **interpolate** for any arbitrary x in between the labelled x_i ...
 - ...in million of dimensions!

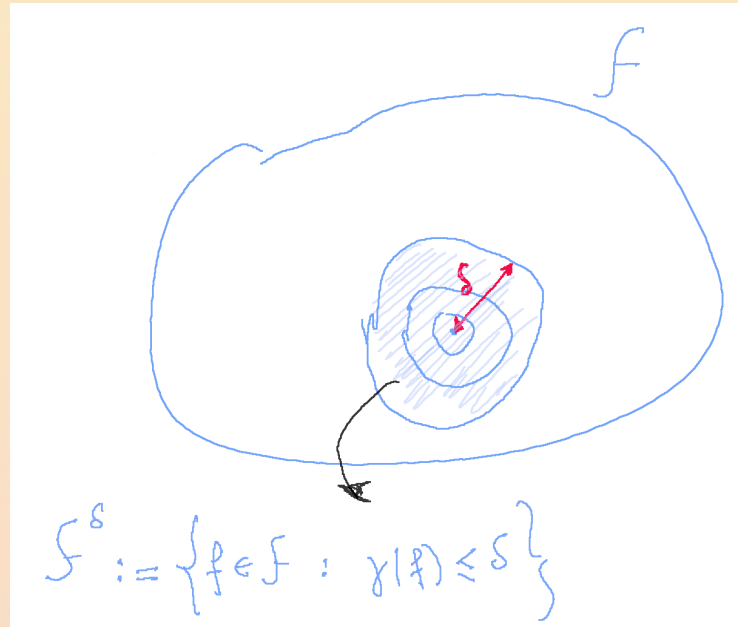


The space of possible solutions

- The space of possible functional solutions is vast: $\mathcal{F} \subseteq \{f : \mathcal{X} \rightarrow \mathbb{R}\}$
(hypothesis class)
- We need a notion of complexity to "organize" the space
- $\gamma(f), f \in \mathcal{F}$: **complexity** of f
 - It can for example be the norm, i.e. we can augment the space \mathcal{F} with the norm

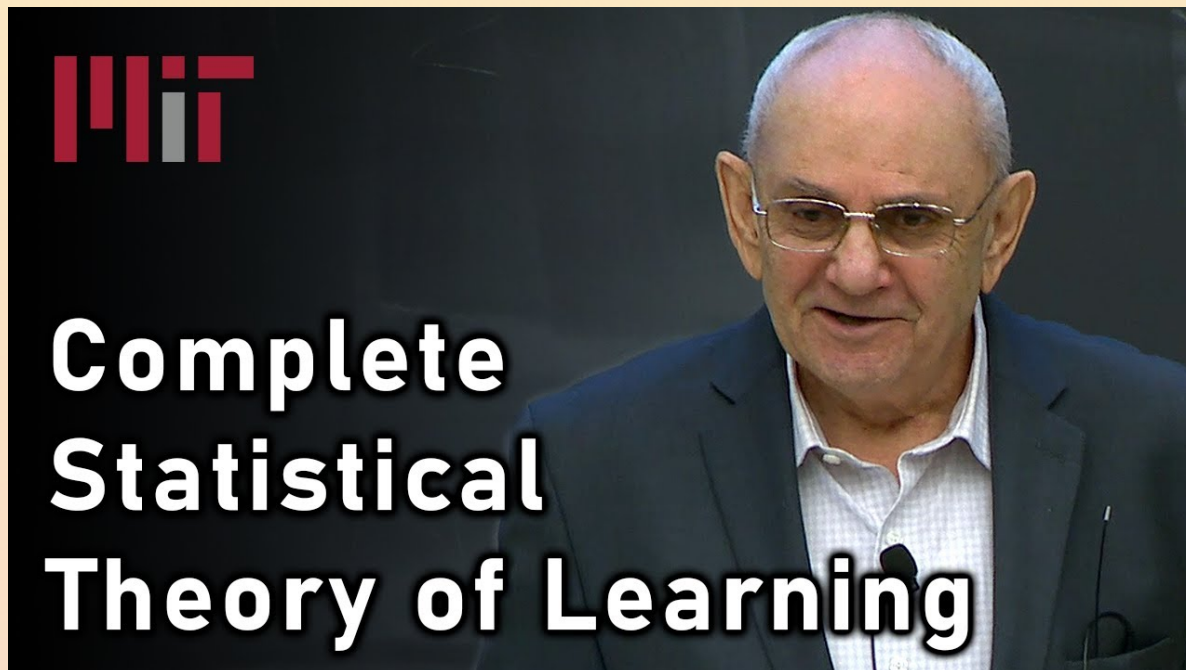
Organizing the space

- When the complexity is defined via the norm, \mathcal{F} is highly organized: Banach space!
 - The simplest function according to the norm criterion is the $\mathbf{0}$ function
 - If we increase the complexity by increasing the norm, we obtain **convex balls**
 $\{f \in \mathcal{F}; \gamma(f) \leq \delta\} =: \mathcal{F}^\delta$
 - Convex minimization is considerably easier than non-convex minimization



Empirical Risk Minimization

- For each element of \mathcal{F} , a measure of how well it's interpolating the data
- Empirical risk: $\hat{L}(f) = \frac{1}{n} \sum_{i=1}^n |f(x_i) - f^*(x_i)|^2$
 - $|\cdot|$ is the empirical loss. If it's the norm, then $\hat{L}(f)$ is the empirical Mean Square Error
 - If you find an analogy with least squares method, it's because for one variable it's exactly that!



Formalizing the minimization of a functional in a given space}

- Constraint form: $\min_{f \in \mathcal{F}^\delta} \hat{L}(f)$.
 - Not trivial
- Penalized form: $\min_{f \in \mathcal{F}} \hat{L}(f) + \lambda \gamma(f)$.
 - More typical
 - λ is the price to pay for more complex solutions. Depends on the complexity measure
- Interpolant form: $\min_{f \in \mathcal{F}} \gamma(f)$ s.t. $\hat{L}(f) = 0 \iff f(x_i) = f^*(x_i) \quad \forall i$
 - In ML, most of the times there is no noise, so $f(x_i)$ is exactly the value we expect there (i.e. we really know that x_i is of a given class, without any uncertainty)

Are the three forms equivalent?

- The interpolant form exploits the no-noise assumption (*"give me the least complex elements in \mathcal{F} that interpolates"*)
- These forms are **not completely equivalent**. The penalized form to be solved requires averaging a full set of penalized forms, so it's not completely equivalent
- There is certainly an implicit correspondence between δ and λ
 - (the larger λ , the smaller δ and viceversa)

The Fundamental Theorem of Machine Learning

- We want to relate the result of the empirical risk minimization (ERM) with the prediction
 - Let's use the constraint form
- Let's assume we have solved the ERM at a precision ϵ (we are ϵ -away from...). We then have $\hat{f} \in \mathcal{F}^\delta$ such that $\hat{L}(\hat{f}) \leq \epsilon + \min_{f \in \mathcal{F}^\delta} \hat{L}(f)$
- How good is \hat{f} at predicting f^* ? In other words, what's the true loss?
 - Can use the triangular inequality

$$L(\hat{f}) - \inf_{f \in \mathcal{F}} L(f) \leq \inf_{f \in \mathcal{F}^\delta} L(f) - \inf_{f \in \mathcal{F}} L(f)$$

$$+ 2 \sup_{f \in \mathcal{F}^\delta} |L(f) - \hat{L}(f)|$$

$$+\epsilon$$

Approximation error

(how appropriate is my measure of complexity)

Statistical error

(having the empirical loss instead of the true loss)

Optimization error

The Error Market

- The minimization is regulated by the parameter δ (the size of the ball in the space of functions)
- Changing δ results in a **tradeoff** between the different errors
 - Very small δ makes the statistical error blow up
- We are better at doing convex optimization (easier to find minimum), but even then the optimization error ϵ will not be negligible
 - ϵ : how much are you willing to spend in resources to minimize $\hat{L}(f)$
 - We kind of control it!
 - If the other errors are smaller than ϵ , then it makes sense to spend resources to decrease it
 - Otherwise, don't bother

The big questions

- **Approximation**: we want to design "good" spaces \mathcal{F} to approximate f^* in high-dimension
 - Rather profound problem, on which we still struggle
- **Optimization**: how to design algorithms to solve the ERM in general
 - We essentially have ONE answer
 - **Gradient Descent!**

The Curse of Dimensionality

- How many samples do we need to estimate f^* , depending on assumptions on its regularity?

The Curse of Dimensionality

- How many samples do we need to estimate f^* , depending on assumptions on its regularity?
 - f^* constant \rightarrow need only 1 sample
 - f^* linear \rightarrow need d samples
- Space of functionals is $\mathcal{F} = \{f : \mathbb{R}^d \rightarrow \mathbb{R}; f(x) = \langle x, \theta \rangle\} \simeq \mathbb{R}^d$ (isomorphic)
 - It's essentially like solving a system of linear equations for the linear form $\langle x, \theta^* \rangle$
 - d equations, d degrees of freedom
- The reason why it's so easy is that linear functions are regular at a global level
 - Knowing the function locally tells us automatically the properties everywhere

Locally linear functions

- Assume f^* locally linear, i.e. f^* is Lipschitz
 - $|f^*(x) - f^*(y)| \leq \beta \|x - y\|$
 - $Lip(f^*) = \inf \{ \beta; |f^*(x) - f^*(y)| \leq \beta \|x - y\| \text{ is true} \}$
 - $Lip(f^*)$ is a measure of smoothness
- Space of functionals that are Lipschitz: $\mathcal{F} = \{ f : \mathbb{R}^d \rightarrow \mathbb{R}; f \text{ is Lipschitz} \}$
- We want a normed space to parameterize complexity, so we convert to a Banach space
 - $\gamma(f) := \max(Lip(f), \|f\|_\infty)$
 - The parameterization of complexity is the Lipschitz constant

Formalization of the prediction problem

- $\forall \epsilon > 0$, find $f \in \mathcal{F}$ such that $\|f - f^*\| \leq \epsilon$ from n i.i.d. samples
 - n : sample complexity, "how many more samples to I need to make the error a given amount of times smaller"
- If f^* is Lipschitz, it can be demonstrated that $n \sim \epsilon^{-d}$
 - Upper bound: approximate f with its value in the closest of the sampled data points, find out expected error $\sim \epsilon^2$, upper bound is exponential
 - Lower bound: maximum discrepancy (the worst case scenario): unless you sample exponential number of data points, knowing $f(x_i)$ for all of them doesn't let you well approximate outside

PAC learning

DEFINITION 3.1 (PAC Learnability) A hypothesis class \mathcal{H} is PAC learnable if there exist a function $m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathbb{N}$ and a learning algorithm with the following property: For every $\epsilon, \delta \in (0, 1)$, for every distribution \mathcal{D} over \mathcal{X} , and for every labeling function $f : \mathcal{X} \rightarrow \{0, 1\}$, if the realizable assumption holds with respect to $\mathcal{H}, \mathcal{D}, f$, then when running the learning algorithm on $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ i.i.d. examples generated by \mathcal{D} and labeled by f , the algorithm returns a hypothesis h such that, with probability of at least $1 - \delta$ (over the choice of the examples), $L_{(\mathcal{D}, f)}(h) \leq \epsilon$.

- ϵ ("approximately correct"): how far from optimality the model is
- δ ("probably"): how likely the model is to meet the accuracy requirement
- $m_{\mathcal{H}}$ determines sample complexity (how many examples to guarantee PAC?)

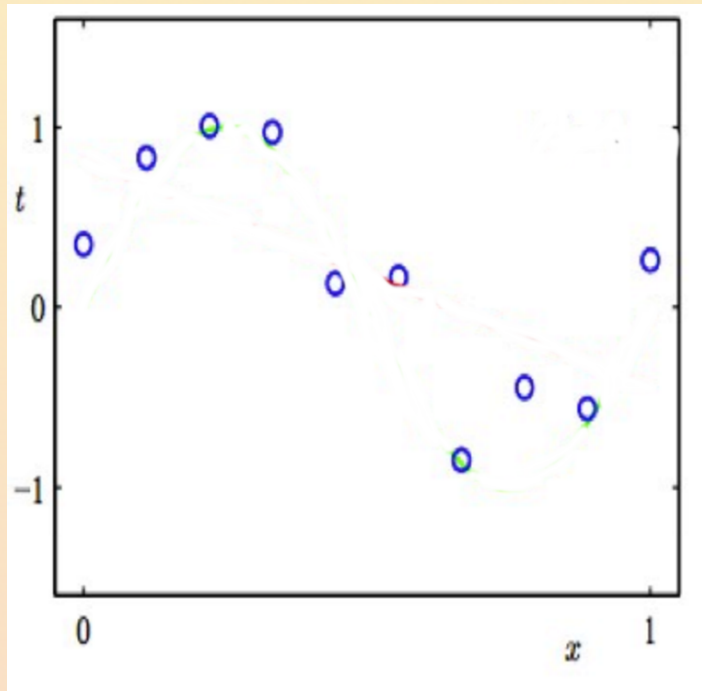
COROLLARY 3.2 *Every finite hypothesis class is PAC learnable with sample complexity*

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{\log(|\mathcal{H}|/\delta)}{\epsilon} \right\rceil.$$

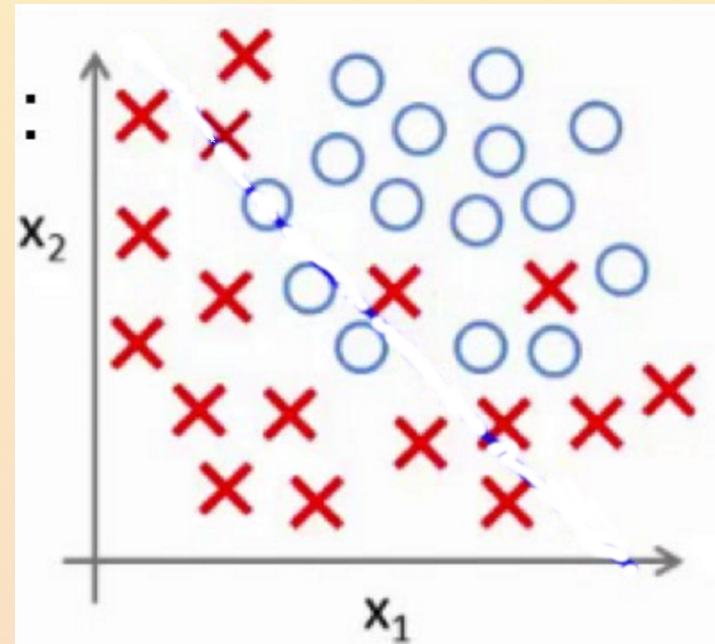
Enough of the math?

What's the best function

To describe the data points?
(regression)

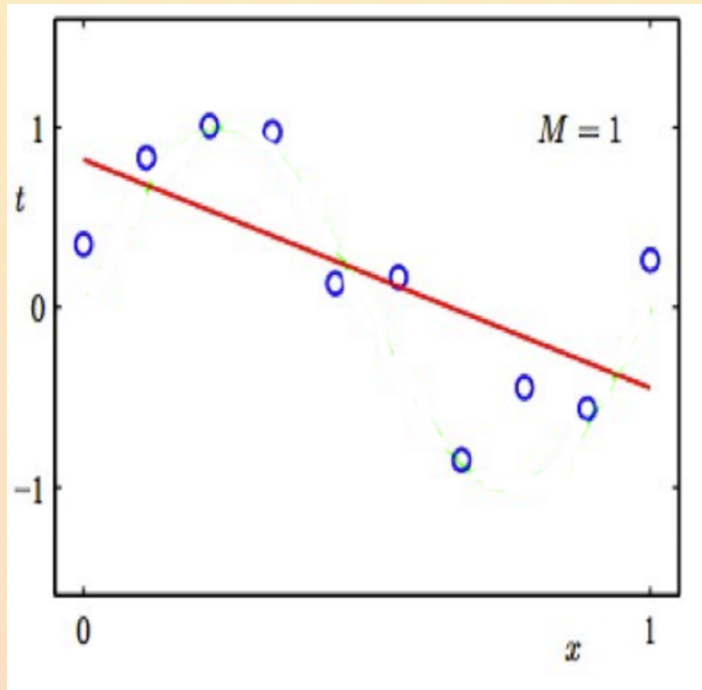


To separate into two classes?
(classification)

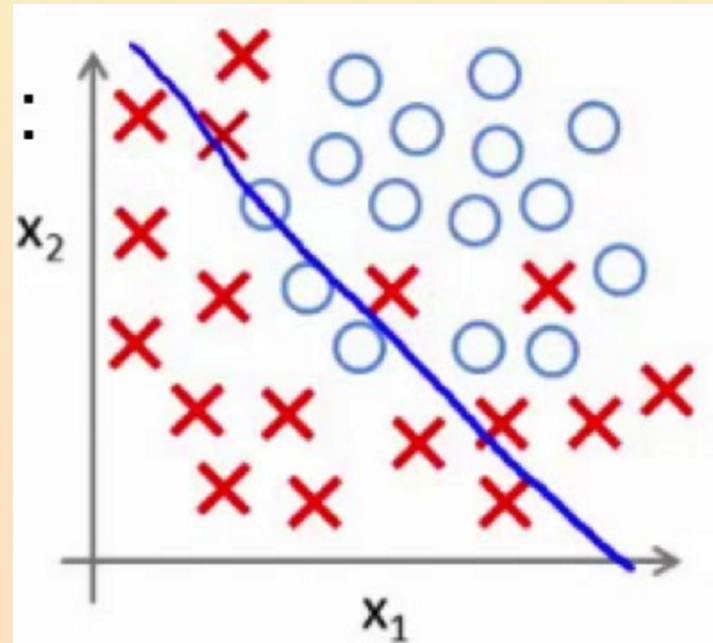


What's the best function

To describe the data points?
(regression)

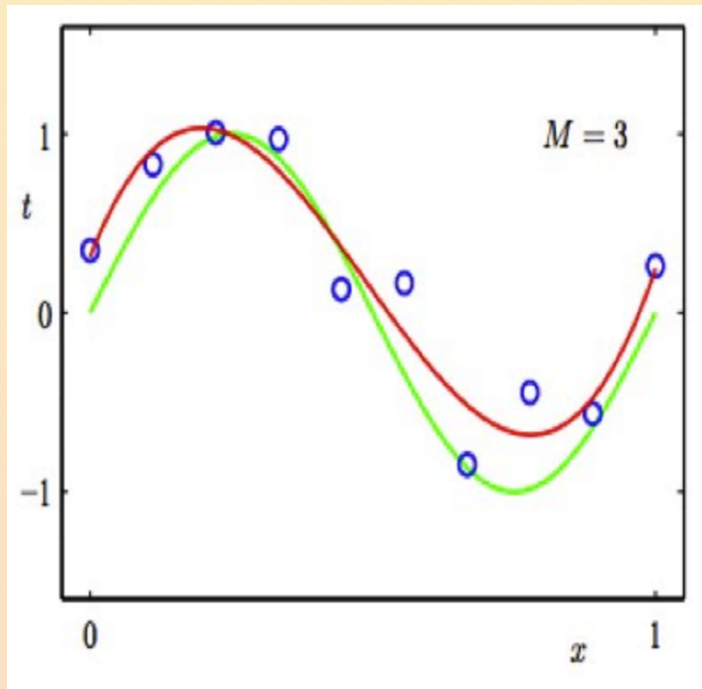


To separate into two classes?
(classification)

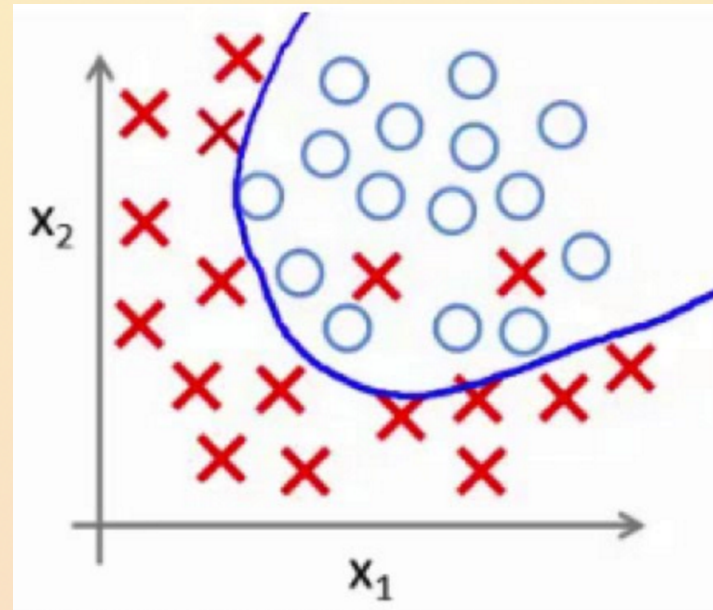


What's the best function

To describe the data points?
(regression)

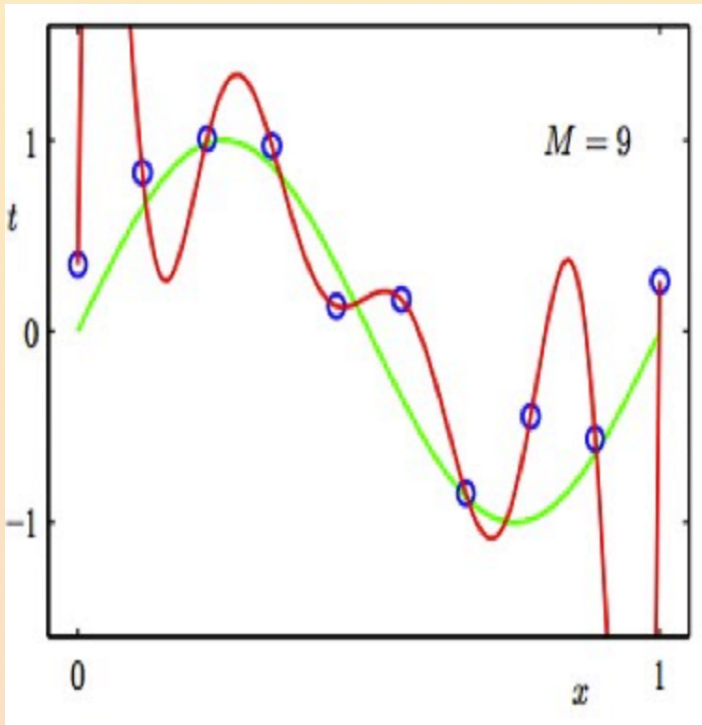


To separate into two classes?
(classification)

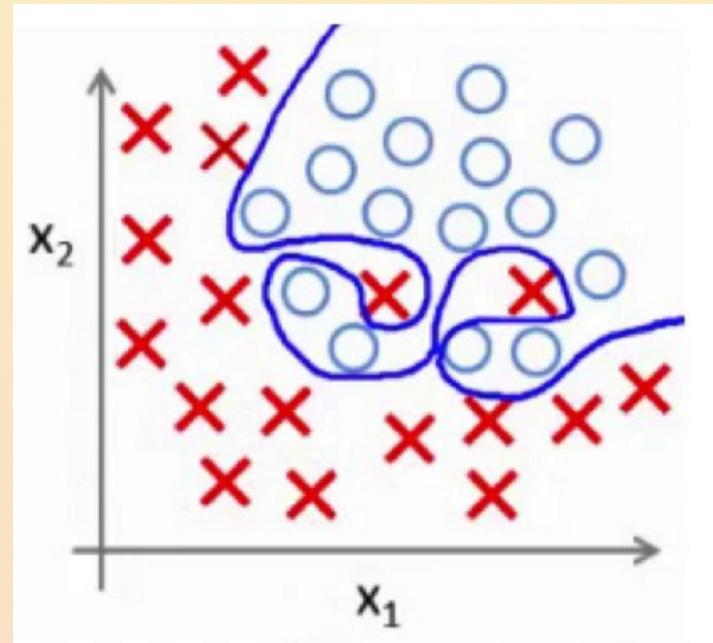


What's the best function

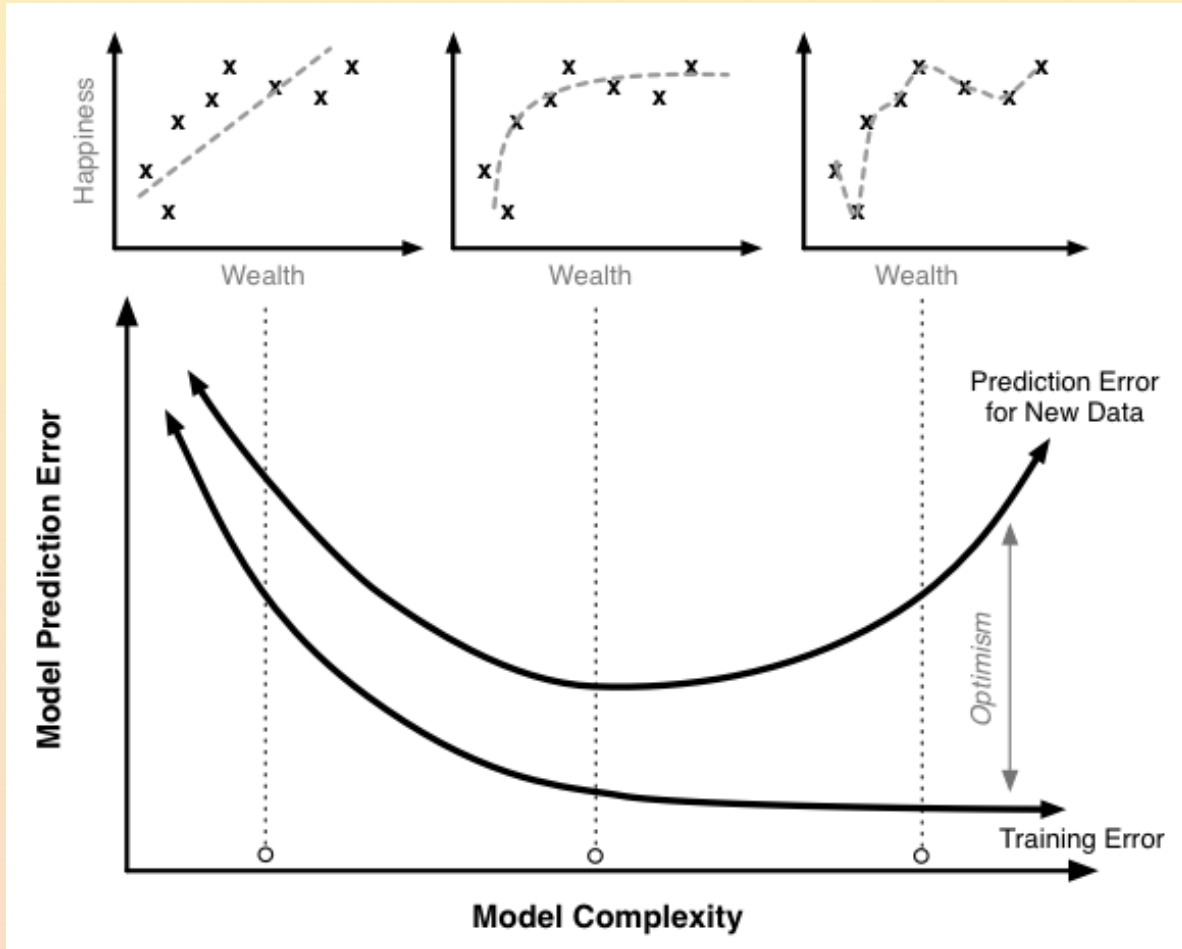
To describe the data points?
(regression)



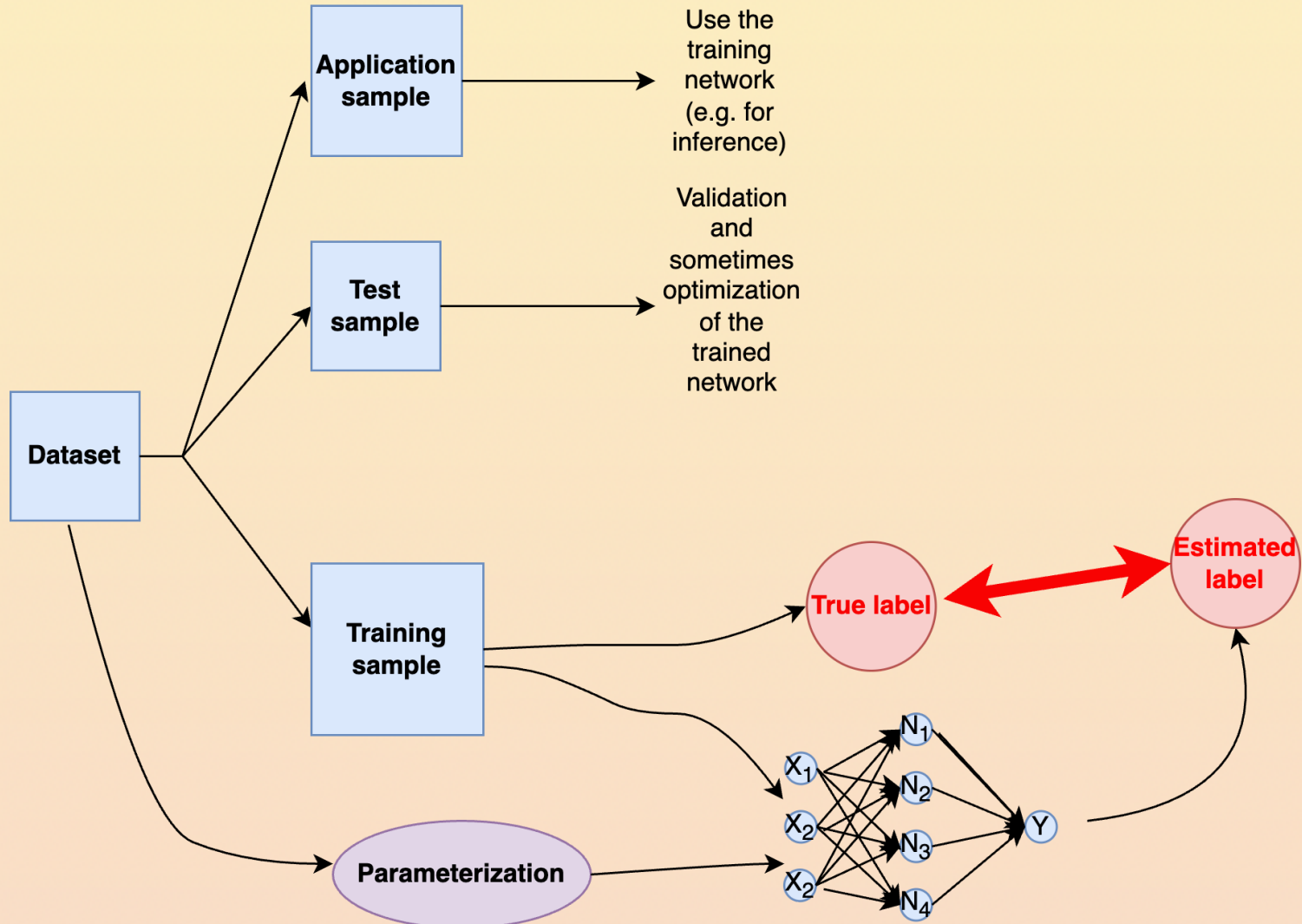
To separate into two classes?
(classification)



Avoid overtraining



Training a model

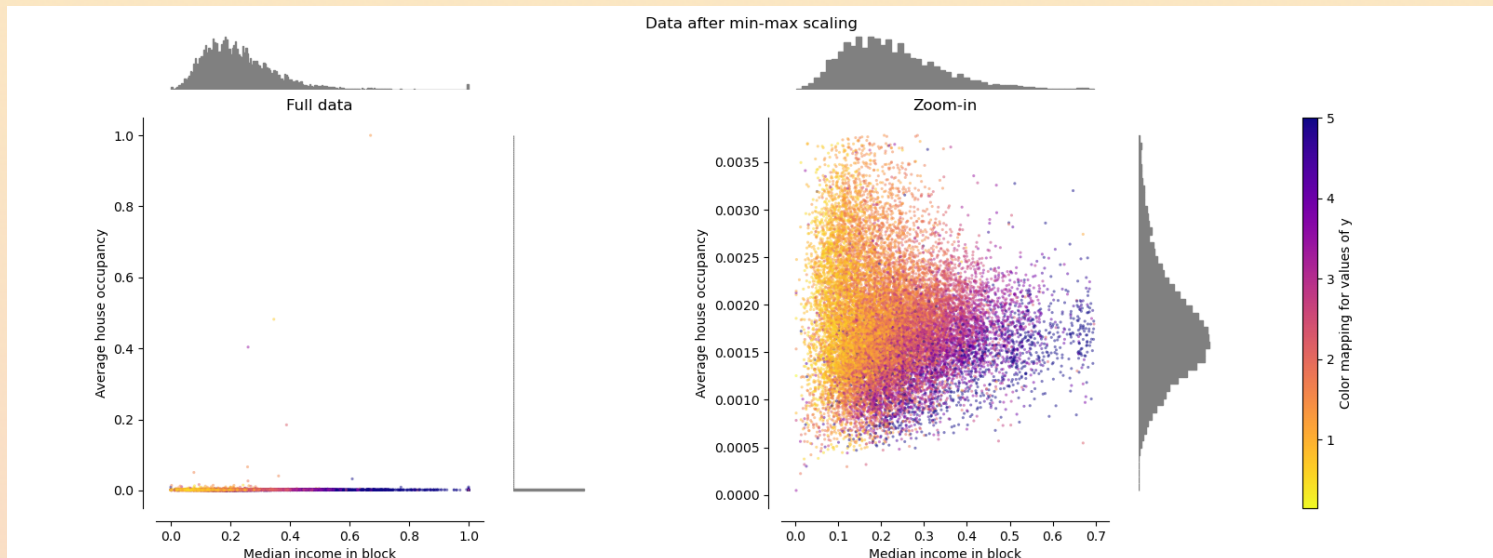


Data preparation: MinMax scaling

$$k = \frac{(x - \min(x))}{\text{abs}(\max(x) - \min(x))}$$

$$z = k * (\text{Max} - \text{Min}) + \text{Min}$$

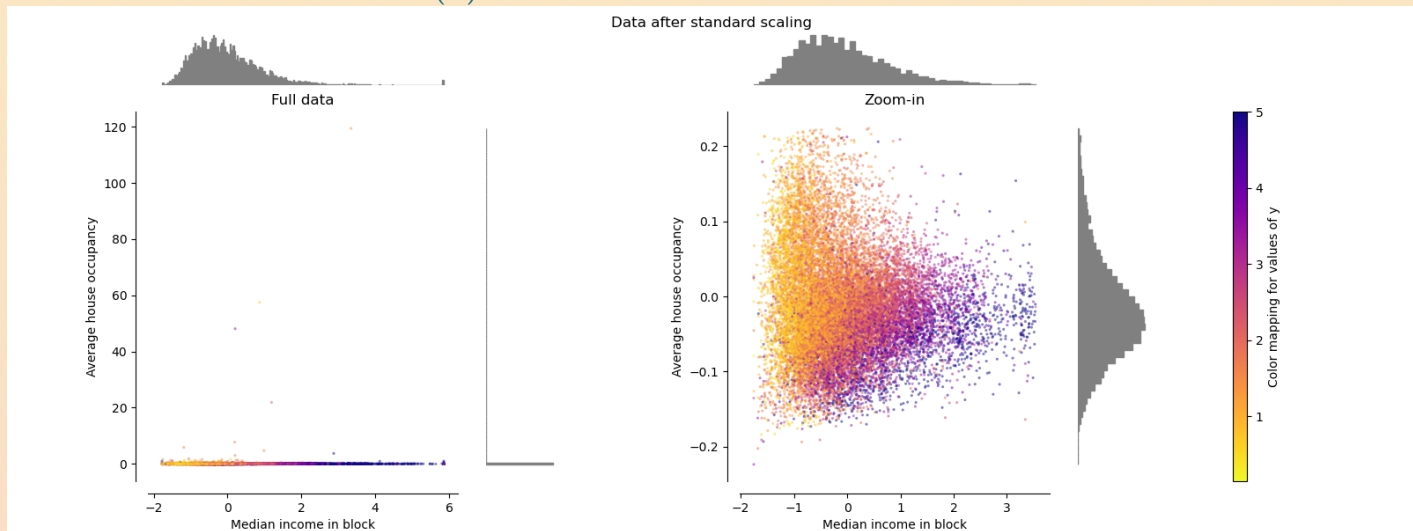
- Very sensitive to outliers (it scales them linearly)
- Scale to different variance and different mean



Data preparation: standardization

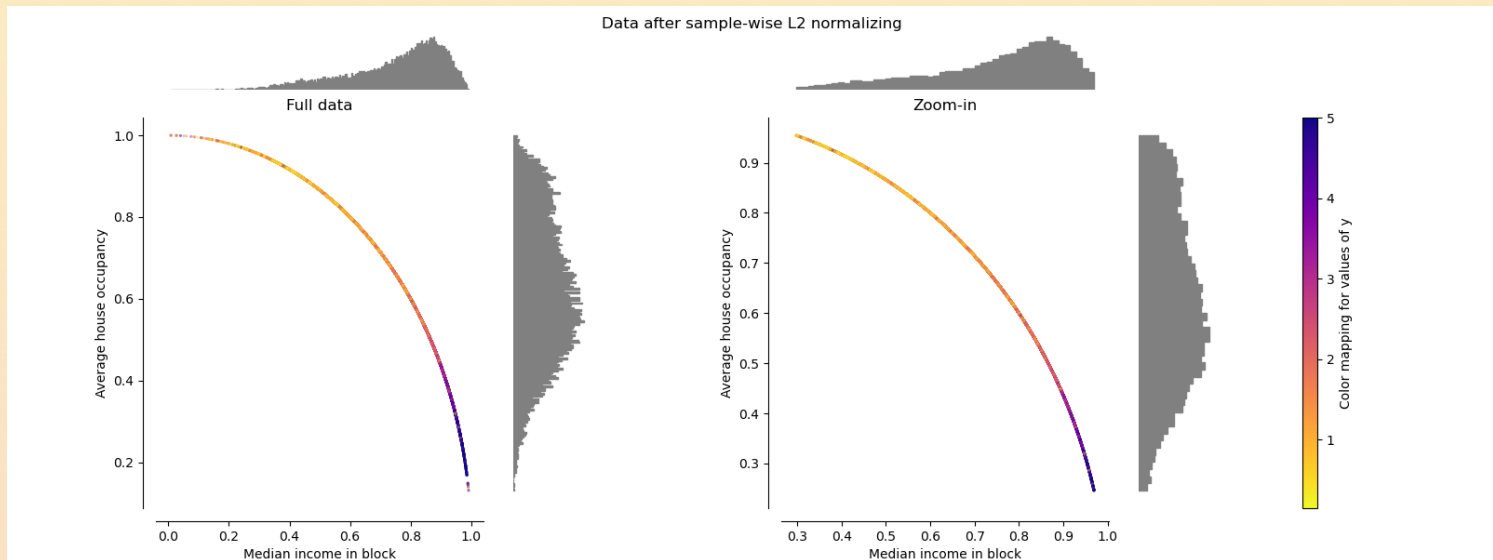
$$z = \frac{(x - \text{mean}(x))}{\text{var}(x)}$$

- Direct comparison between weights assigned to different features
- Easier, more effective numerical minimization, but still sensitive to outliers
- Scale to same variance and same mean (can also scale to same variance but different mean $z = \frac{x}{\text{var}(x)}$)



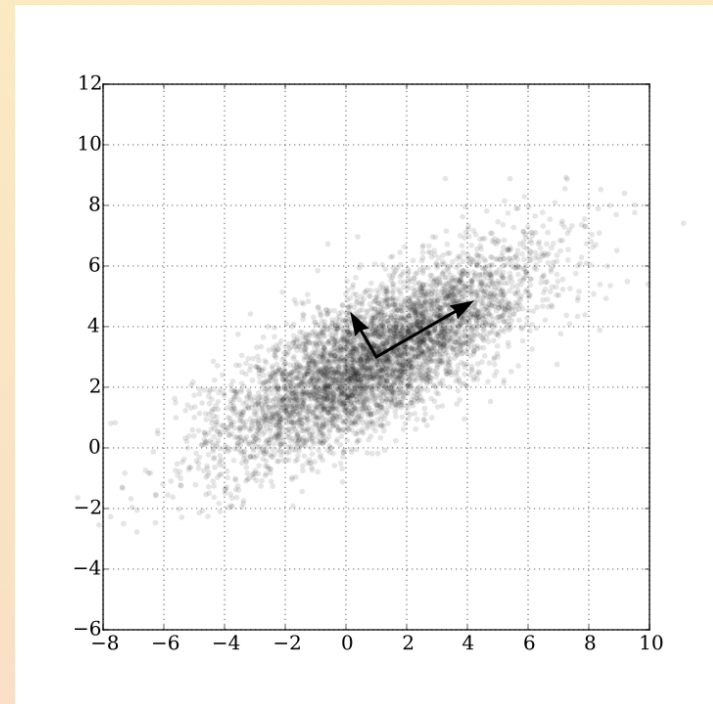
Data preparation: normalisation

- Normalize each data point to have unit norm
- Useful if using dot-product or other kernels to quantify similarity



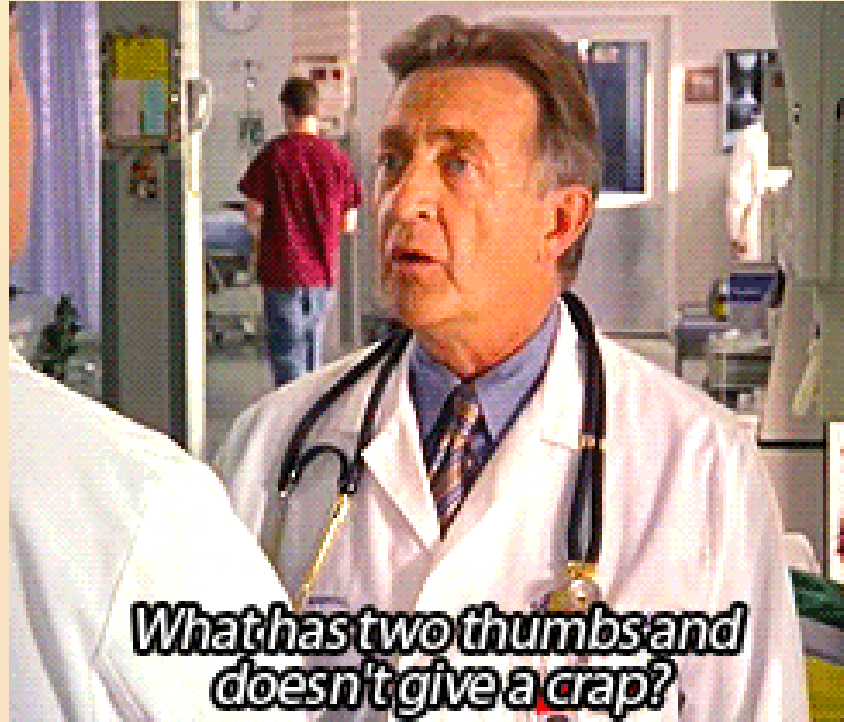
Data preparation: PCA

- Principal Component Analysis
 - Find iteratively the direction (linear combination of features) explaining the most variance
- Principal components are the **eigenvectors of the data covariance matrix**
 - Can be found by Singular Value Decomposition (SVD)
- Somehow analogous to finding axes of ellipsoid
 - Features with different units → arbitrariness (scale them first)
- Can retain a few dimensions: dimensionality reduction
 - Drop directions least explaining the variance



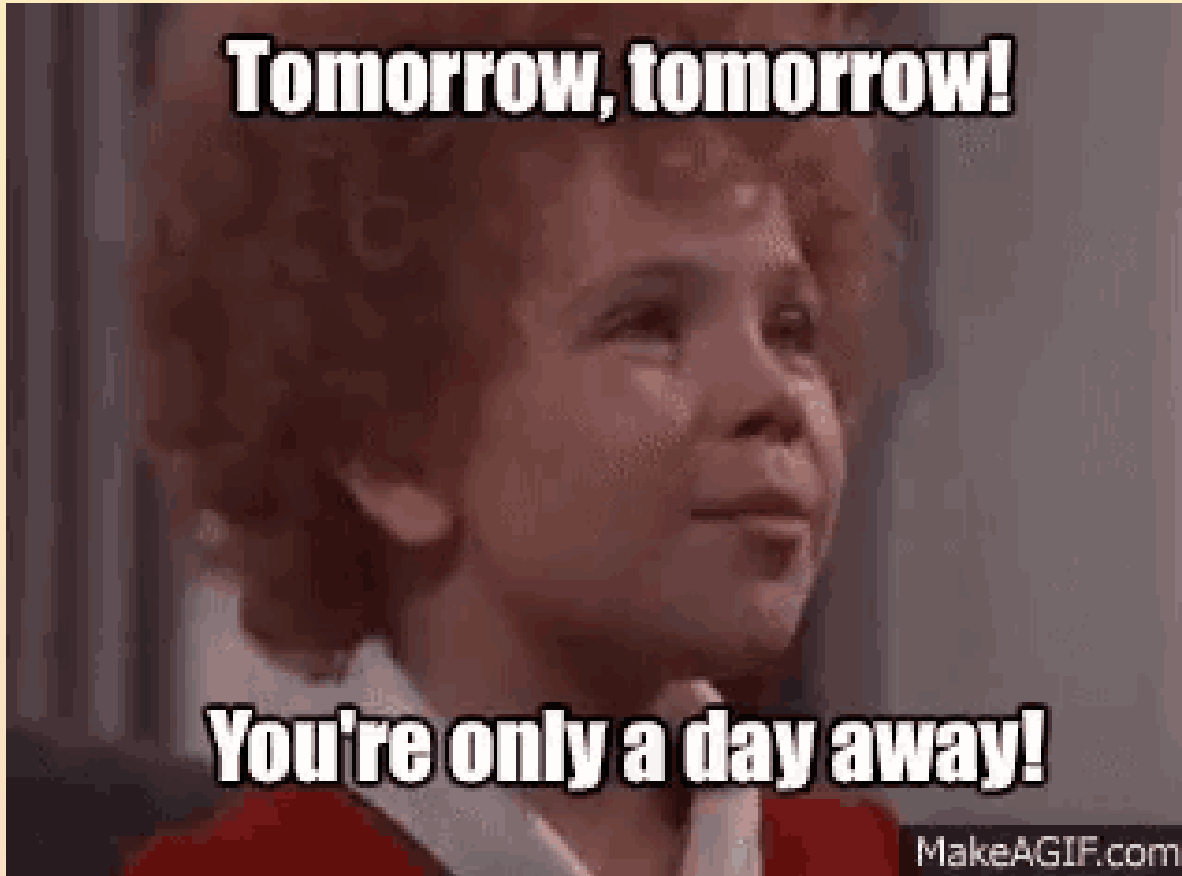
Data preparation: missing values

- LHC data are of **extremely good quality** (unlike e.g. medicine, social sciences)

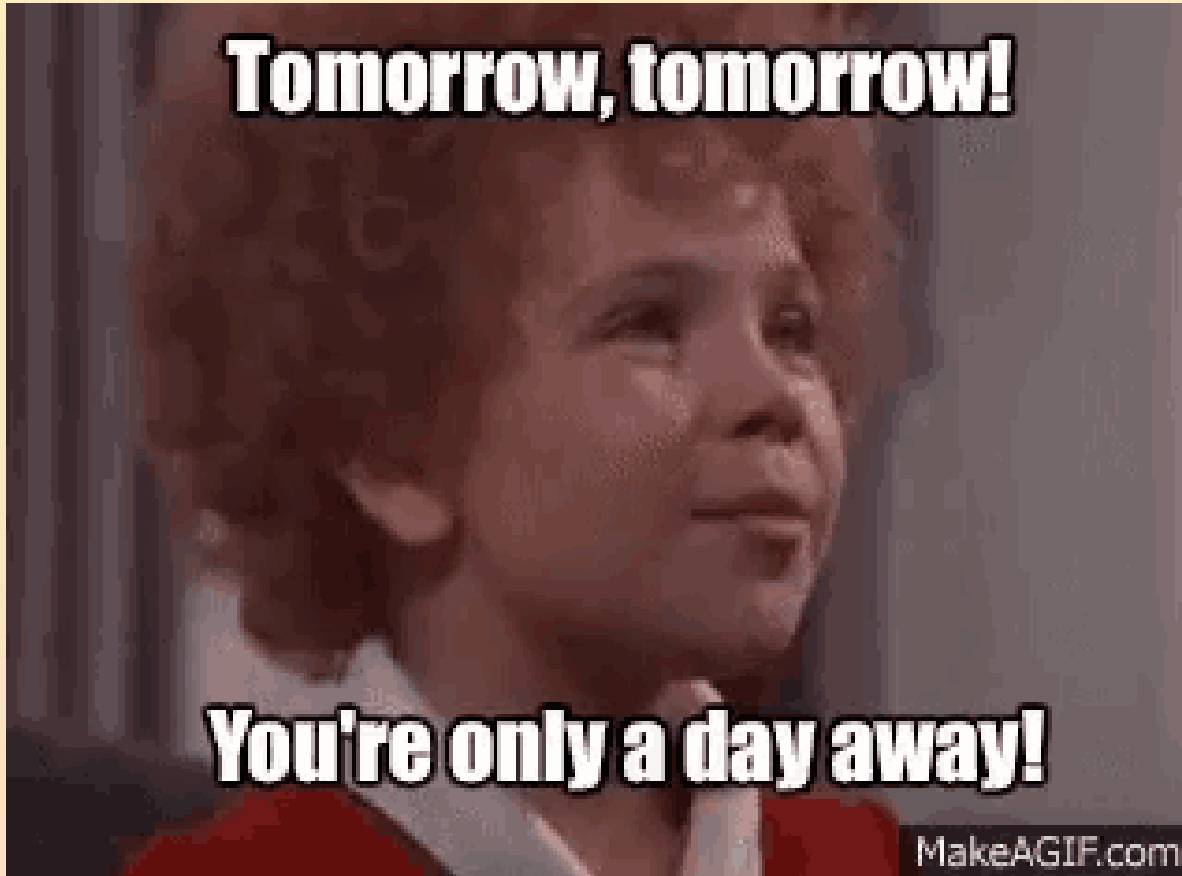


- Use proxy feature (**pT of the two b jets** → **pT of the two jets with highest b-tagging discriminator**, in a region without b jets)
- Use average over other data points (**pT of the third jet** → ****mean of the third jet pT for data points that have it, if some data points don't have three jets**)

Data preparation: feature encoding

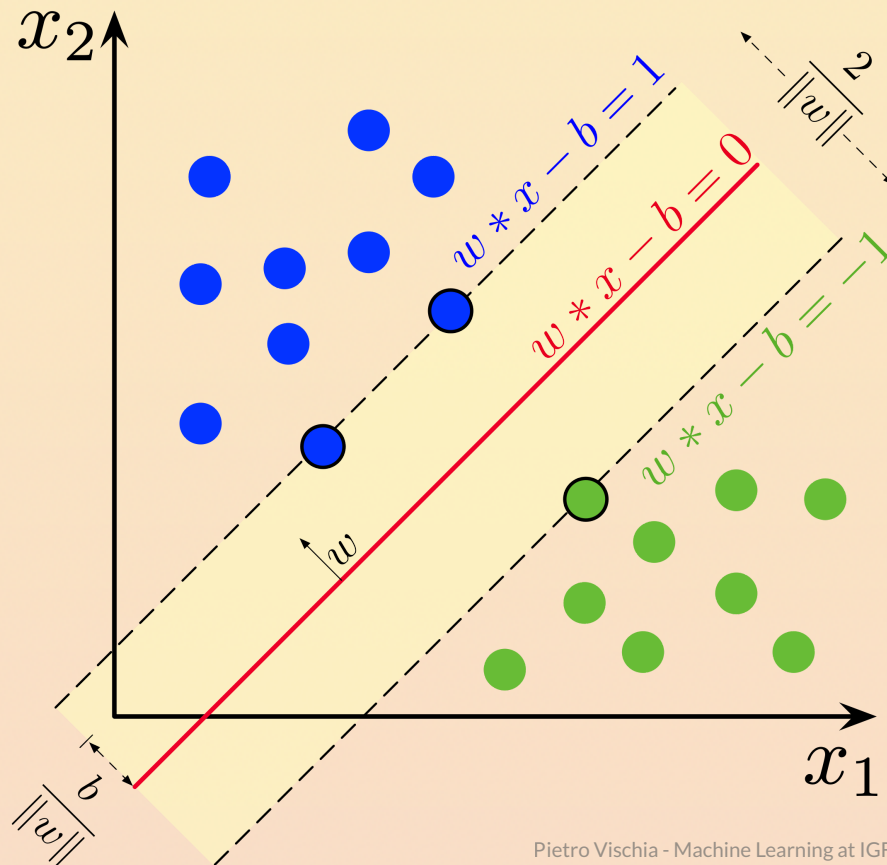


Class weights vs event weights



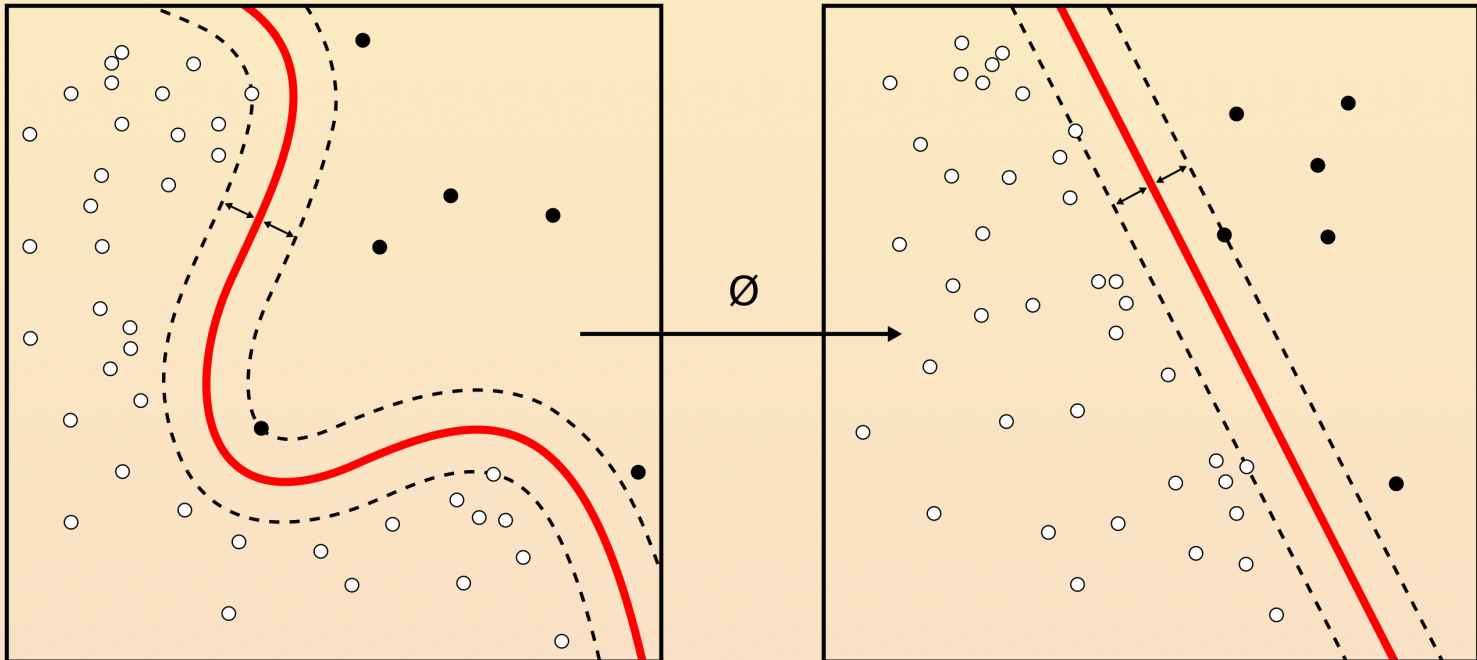
Support Vector Machines (Vapnik)

- Similar to a "cut-based" analysis, but sophisticated strategy for optimal class separation
- Minimize empirical classification error + maximize geometric margin



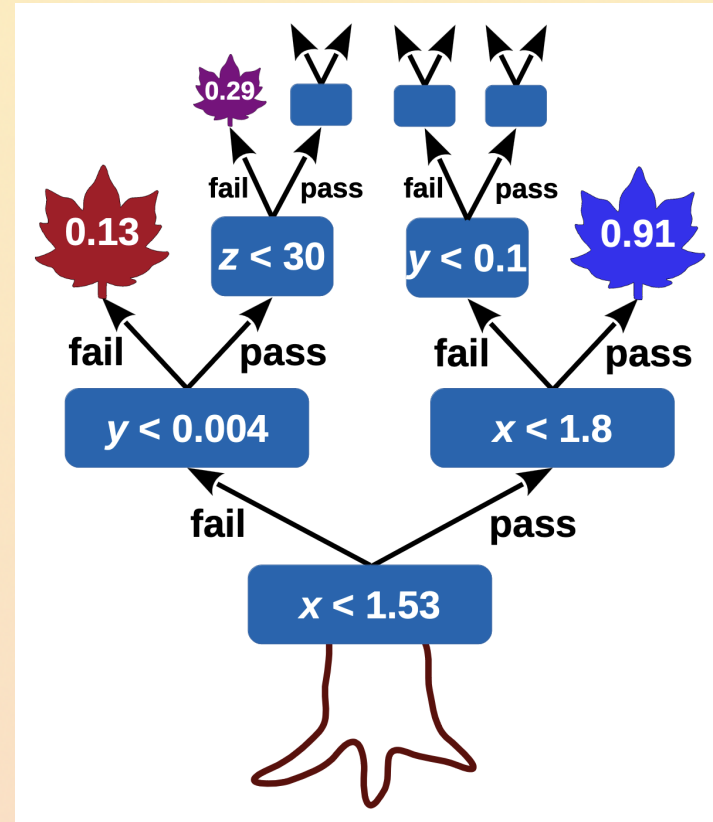
Support vector Machines

- Nonlinearity by kernel trick (deform the metric of the space until separation is linear)



Decision Trees

- "Cut-based" analysis on steroids
- Ordering is key
 - 1) check stopping criterion
 - 2) sort according to each feature
 - 3) compute all separations
 - 4) if best separation improves, split
 - 5) back to 1



Decision Trees: hyperparameters

- Class balance: normalize classes $\sum_i w_i = \sum_j w_j, \forall (i, j)$ at root node
 - In practice, early splitting provides balance anyway
- Impurity $i(t)$, and stopping criterion
 - Minimum leaf (end node) size (maximum error $\sqrt{N_{min}}$ ensures significance of purity estimate)
 - Perfect separation
 - Insufficient improvement
 - Maximum tree depth (purely computational requirement)

Decision Trees: splitting

- Impurity decrease: $\Delta i(S, t) = i(t) - p_P i(t_P) - p_F i(t_F)$
- Optimization problem: $\Delta i(S^*, t) = \max_{S \in \text{splits}} \Delta i(S, t)$

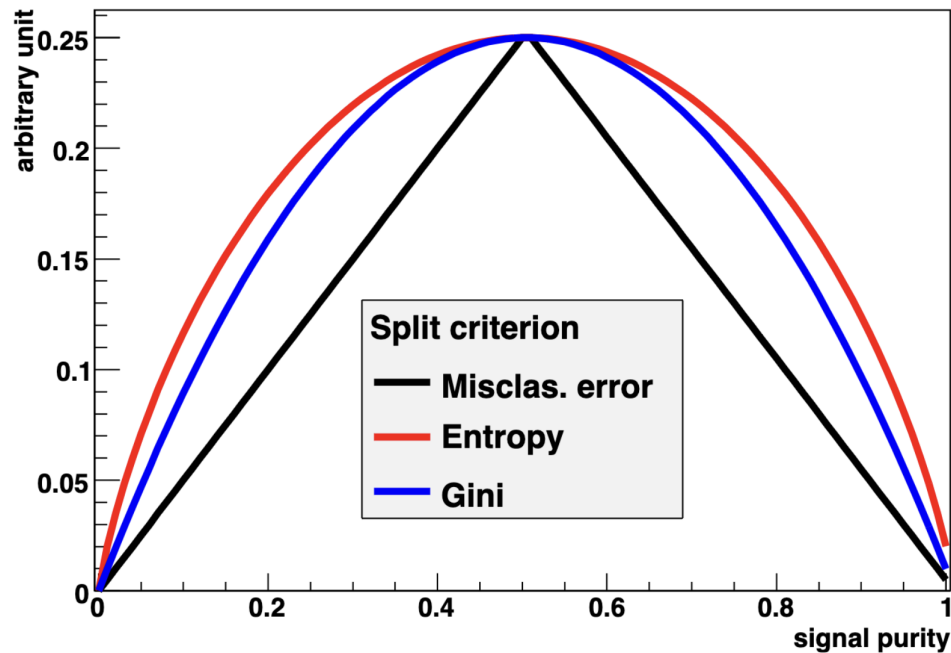
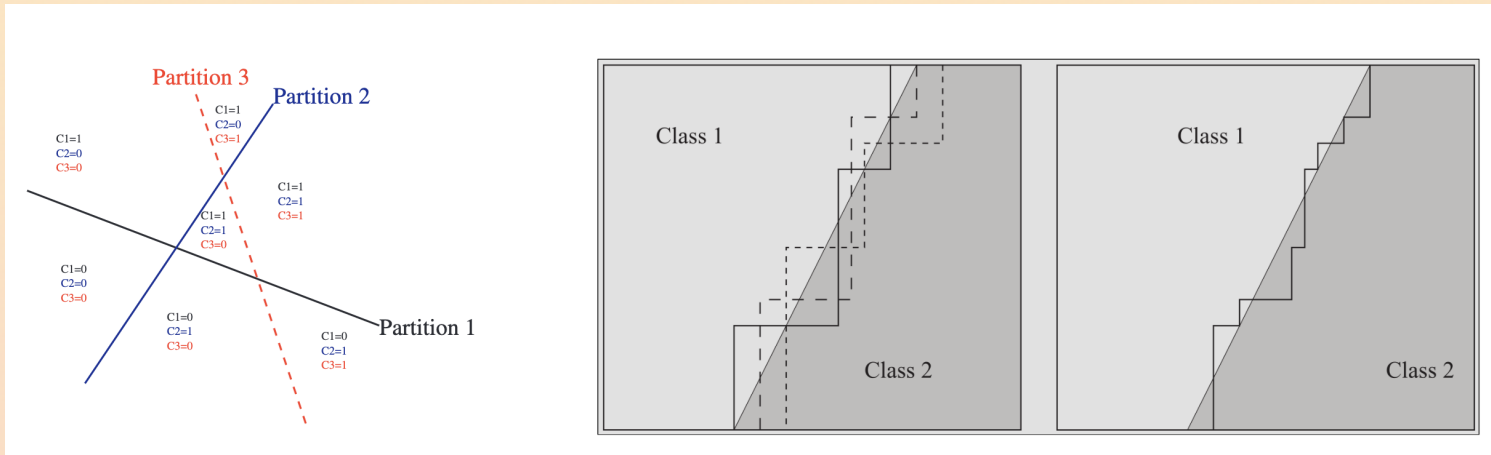


Fig. 5. Impurity measures as a function of signal purity.

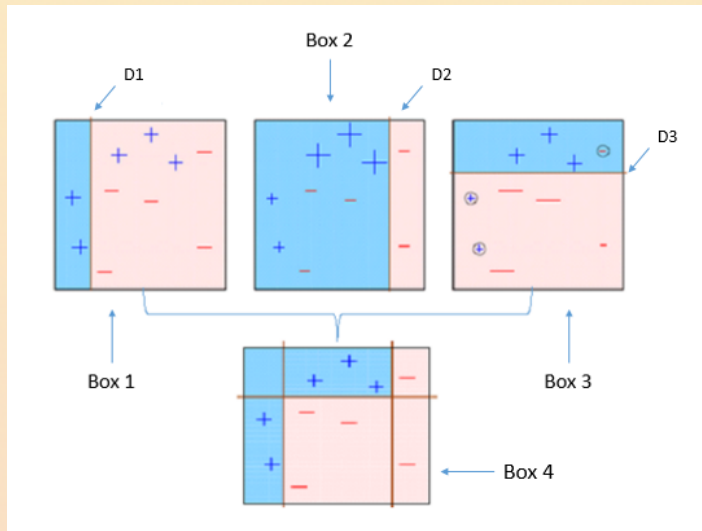
Decision Trees: limitations

- Sensitivity to training sample
 - Decision trees are not robust: high variance for small changes in training sample
- Tree depth results in higher statistical uncertainty in the split purity
 - Can be mitigated by [pruning](#)
- Ensemble learning fixes all of this
 - Richer description when intersecting partitions
 - More accurate description when averaging partitions

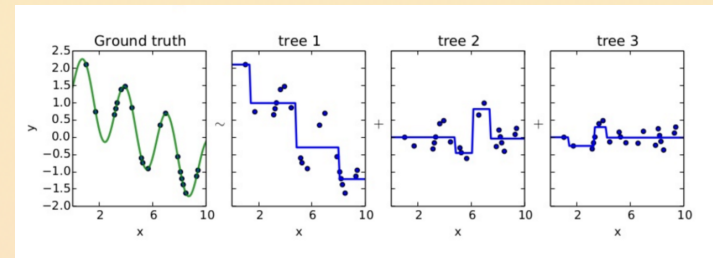


Boosted decision trees

- Ada(ptive) Boost
 - Increase at each iteration the importance of events incorrectly classified in the previous iteration



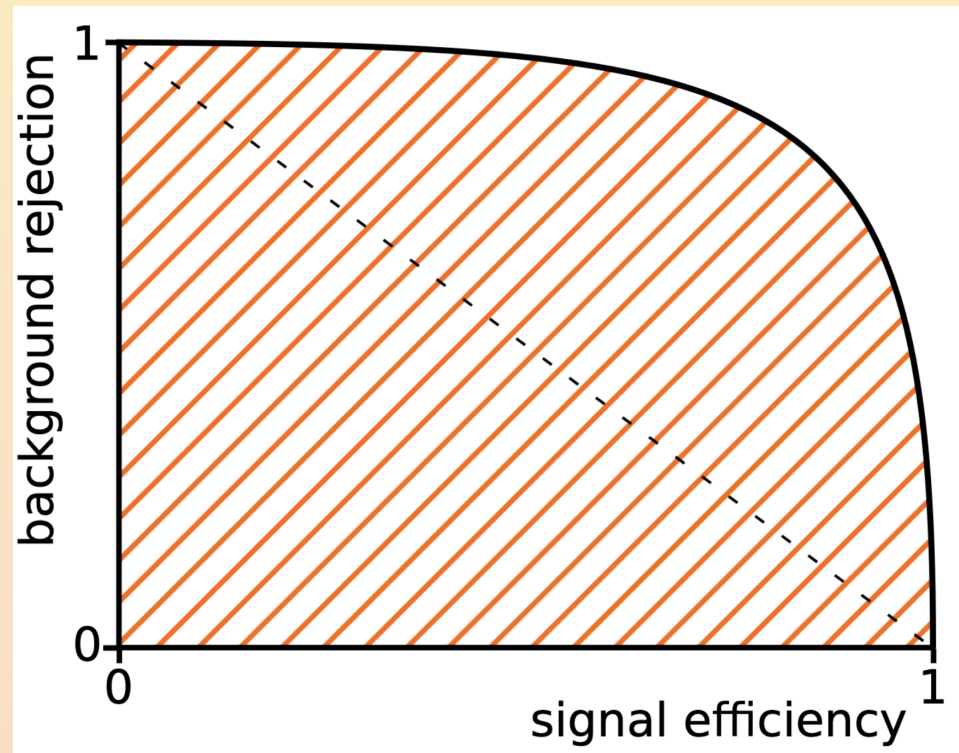
- Gradient Boost
 - fit the new predictor to the residual errors of the previous one



- Bagging: training trees on bootstrap replicas, average all trees
- Random forest: bagging + pick only a random subset of features at each step

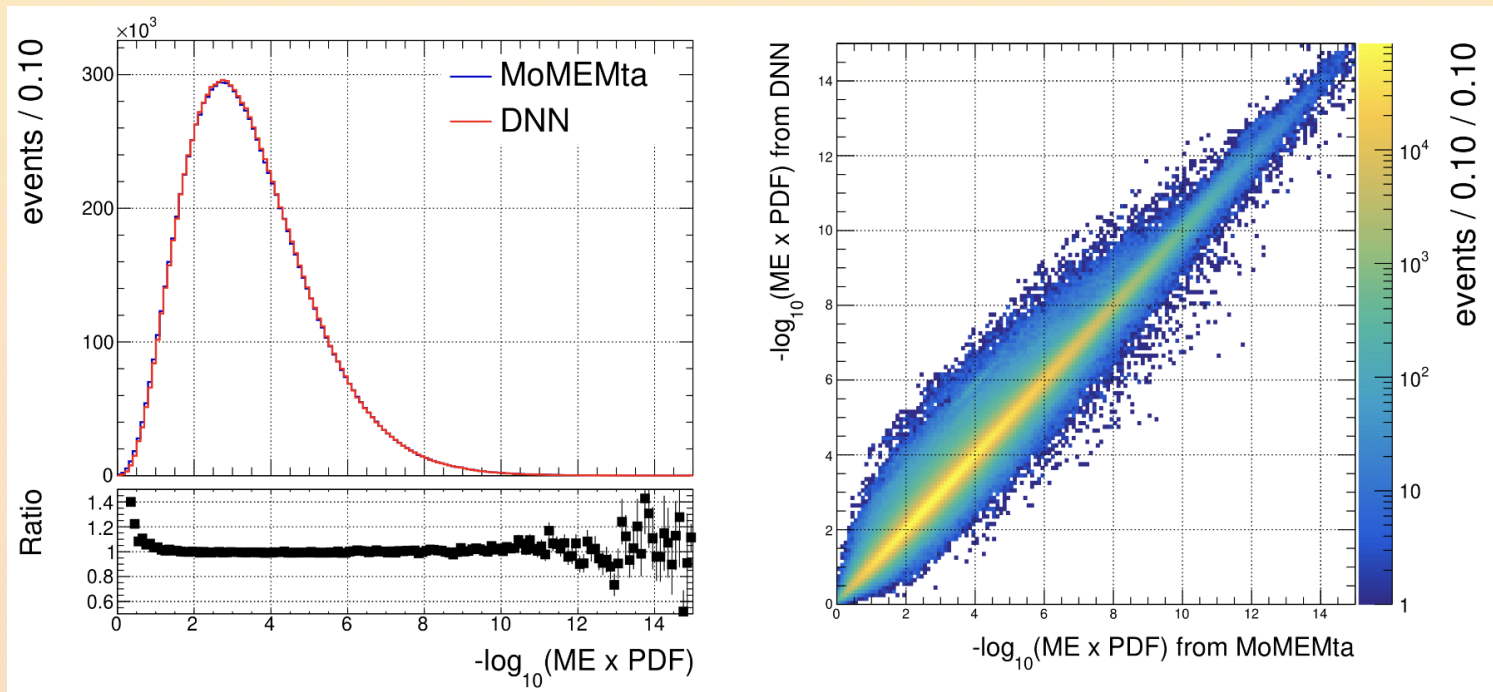
Performance (classification)

- Receiver operating characteristic (ROC) curve
 - Area $\in [0, 1]$, higher values are better
- For two classes (signal, background):
 - Signal efficiency
 - Background efficiency (rejection := 1-efficiency)
- Multiple classes: pairwise, one-vs-all



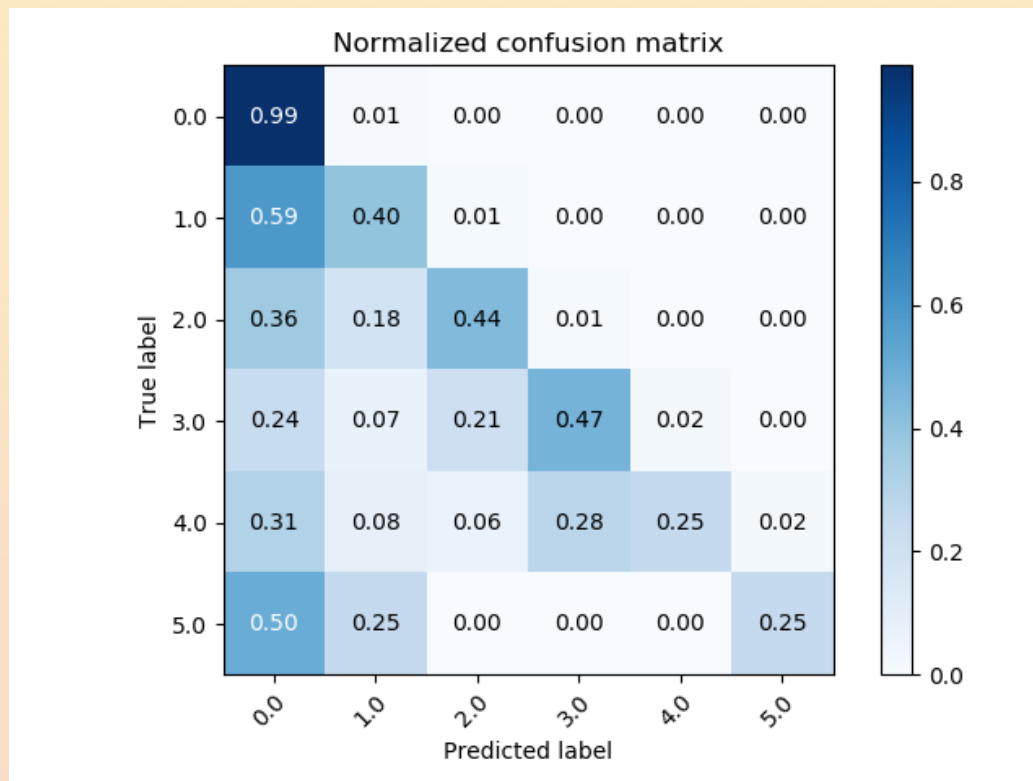
Two-dimensional scatter plot

- For regression problems
- Can compute linear pearson coefficient as an estimate of linearity
 - Formulas exist also for weighted events



Confusion matrix

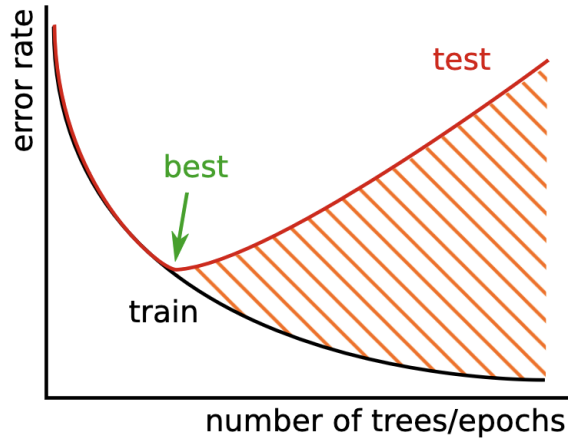
- For classification, can also use it to discretize regression (e.g. in the case of histograms)
- Note the normalization (each true label row sums up to 1)



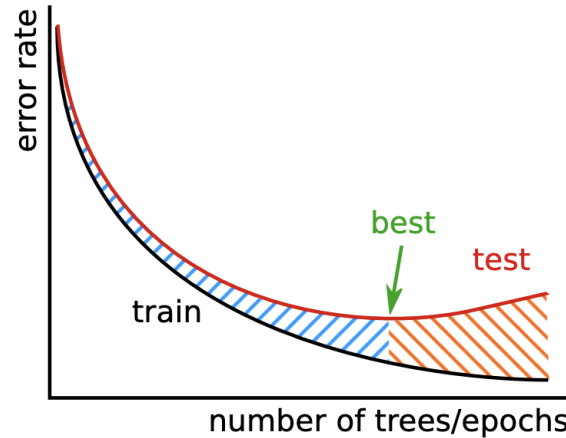
Error rates (for reference)

		CONDITION determined by "Gold Standard"			
TOTAL POPULATION		CONDITION POS	CONDITION NEG	PREVALENCE $\frac{\text{CONDITION POS}}{\text{TOTAL POPULATION}}$	
TEST OUT-COME	TEST POS	True Pos TP	<i>Type I Error</i> False Pos FP	<i>Precision</i> Pos Predictive Value $\text{PPV} = \frac{\text{TP}}{\text{TEST P}}$	False Discovery Rate $\text{FDR} = \frac{\text{FP}}{\text{TEST P}}$
	TEST NEG	<i>Type II Error</i> False Neg FN	True Neg TN	False Omission Rate $\text{FOR} = \frac{\text{FN}}{\text{TEST N}}$	Neg Predictive Value $\text{NPV} = \frac{\text{TN}}{\text{TEST N}}$
ACCURACY ACC $\text{ACC} = \frac{\text{TP} + \text{TN}}{\text{TOT POP}}$		<i>Sensitivity (SN), Recall</i> Total Pos Rate TPR $\text{TPR} = \frac{\text{TP}}{\text{CONDITION POS}}$	<i>Fall-Out</i> False Pos Rate FPR $\text{FPR} = \frac{\text{FP}}{\text{CONDITION NEG}}$	Pos Likelihood Ratio LR + $\text{LR} + = \frac{\text{TPR}}{\text{FPR}}$	Diagnostic Odds Ratio DOR $\text{DOR} = \frac{\text{LR} +}{\text{LR} -}$
		<i>Miss Rate</i> False Neg Rate FNR $\text{FNR} = \frac{\text{FN}}{\text{CONDITION POS}}$	<i>Specificity (SPC)</i> True Neg Rate TNR $\text{TNR} = \frac{\text{TN}}{\text{CONDITION NEG}}$	Neg Likelihood Ratio LR - $\text{LR} - = \frac{\text{TNR}}{\text{FNR}}$	

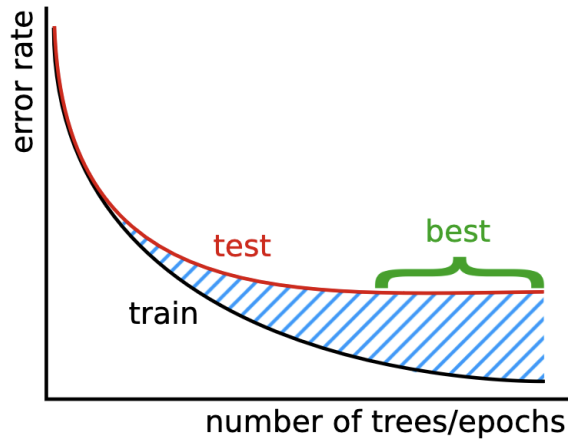
Model complexity



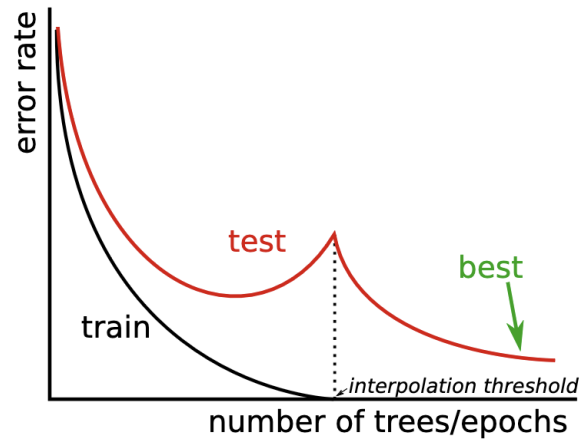
(a)



(b)



(c)

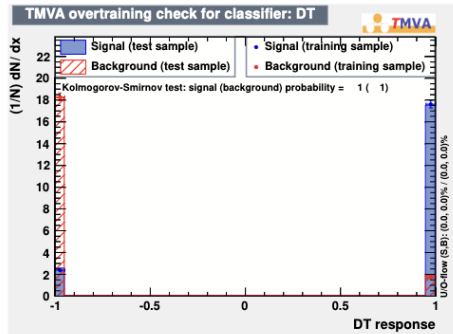


(d)

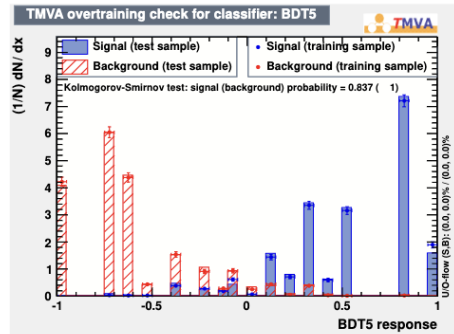
Overtraining check

- KS test done mostly for BDTs. For neural networks, much handier ways

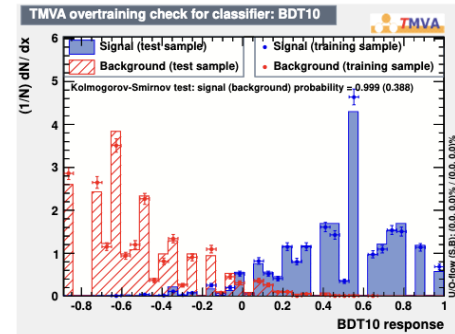
Do the training and test output distributions come from the same underlying p.d.f.?



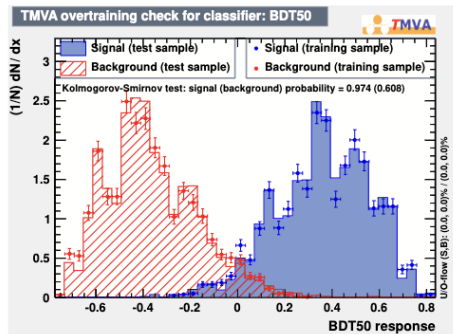
(a) Single decision tree



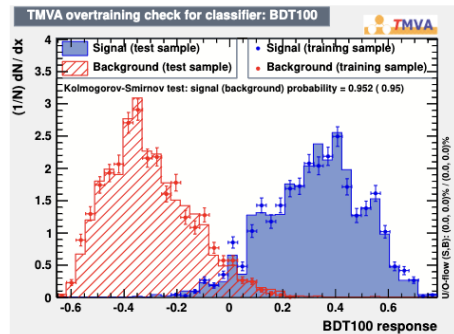
(b) 5 trees



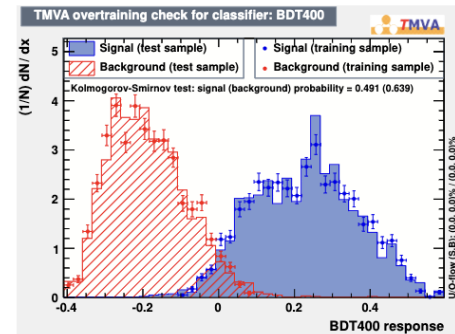
(c) 10 trees



(d) 50 trees



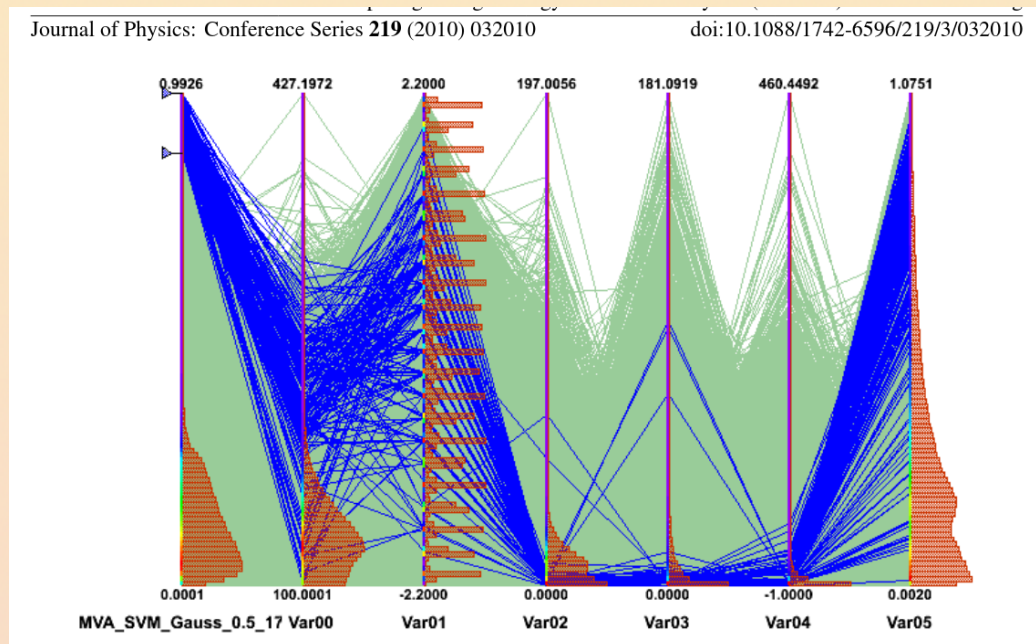
(e) 100 trees



(f) 400 trees

Interpretability, explainability

- **Permutation Importance**: the decrease in a model score when a single feature value is randomly shuffled (see [scikit-learn documentation](#)) (akin to impacts for profile likelihood fits)
- **Shapley Values**: based on game theory (will use them this afternoon)
- **Correlation-based**: e.g. parallel coordinates in TMVA: look where each variable is mapped to/correlated with

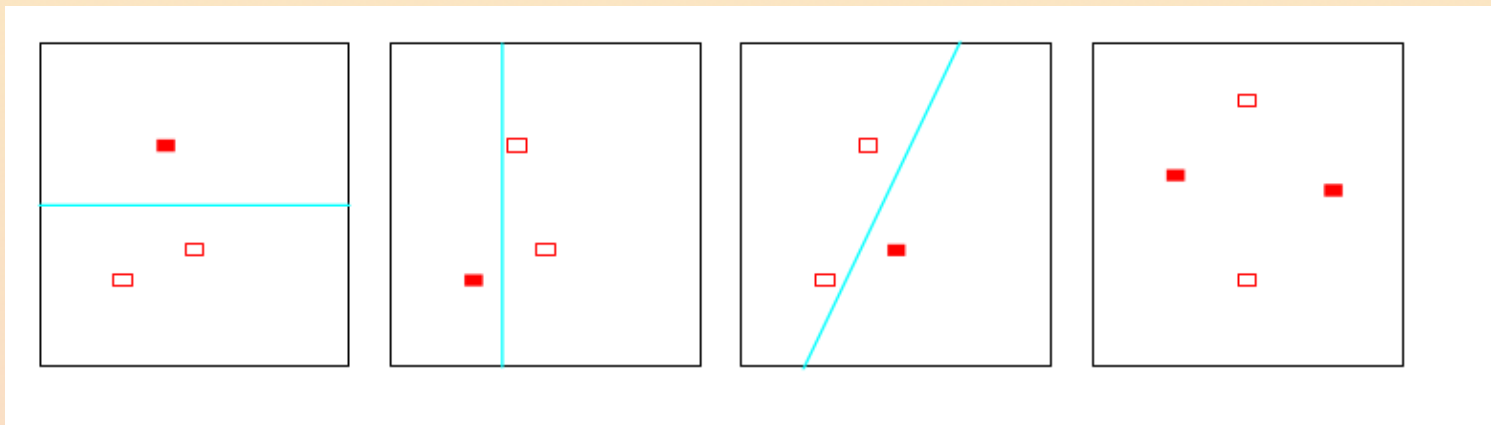


Model assessment by comparing models

- **Bayesian Information Criterion:** $BIC = n_{free\ params} \ln(n_{data}) / 2 \ln(\hat{L})$
 - Parameter θ predicted by two models M_0 and M_1 : $P(\theta|\vec{x}, M) = \frac{P(\vec{x}|\theta, M)P(\theta|M)}{P(\vec{x}|M)}$
 - Apply Bayes theorem to Bayesian evidence (Model likelihood): $P(\vec{x}|M) = \int P(\vec{x}|\theta, M)P(\theta|M)d\theta$
 - Posterior odds: $\frac{P(M_0|\vec{x})}{P(M_1|\vec{x})} = \frac{P(\vec{x}|M_0)\pi(M_0)}{P(\vec{x}|M_1)\pi(M_1)}$
 - Can rewrite posteriors in terms of BIC, equivalent
- **Minimum Description Length (MDL):** Kolmogorov complexity (length of minimum program needed to describe the data)
 - *for i = 1 to 2500; do print'0001'; halt*
 - *print'101001010100010111001000010000101110011100001010100101...'; halt*

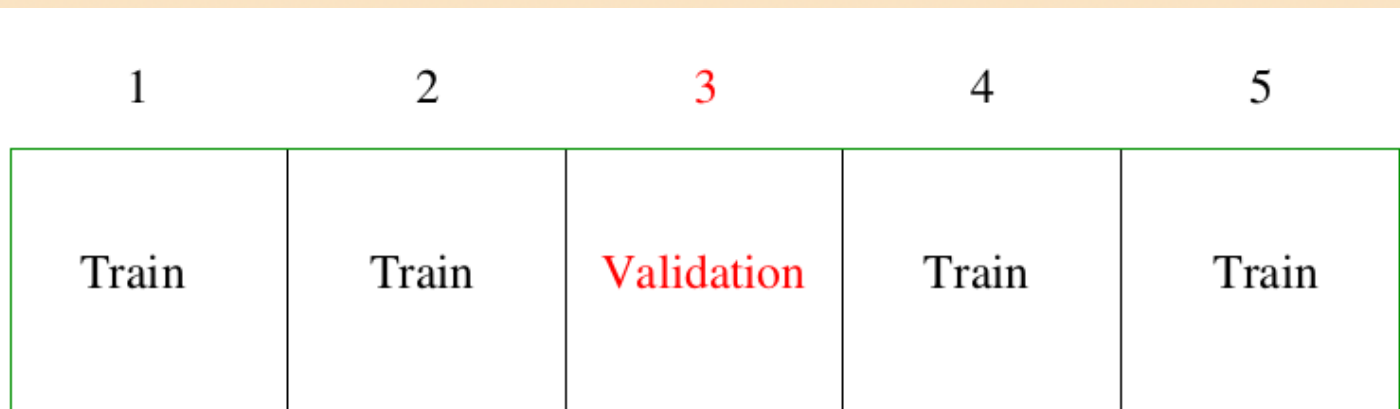
Model assessment by comparing models

- **Structural Risk Minimization:** complexity as Vapnik-Chervonkensis class (largest number of shattered points)
 - Build a nested sequence of models with increasing VC complexity h
 - Write a probabilistic upper bound for the regression error: $err \leq f(h/N)$
 - Choose model with smallest value of the upper bound



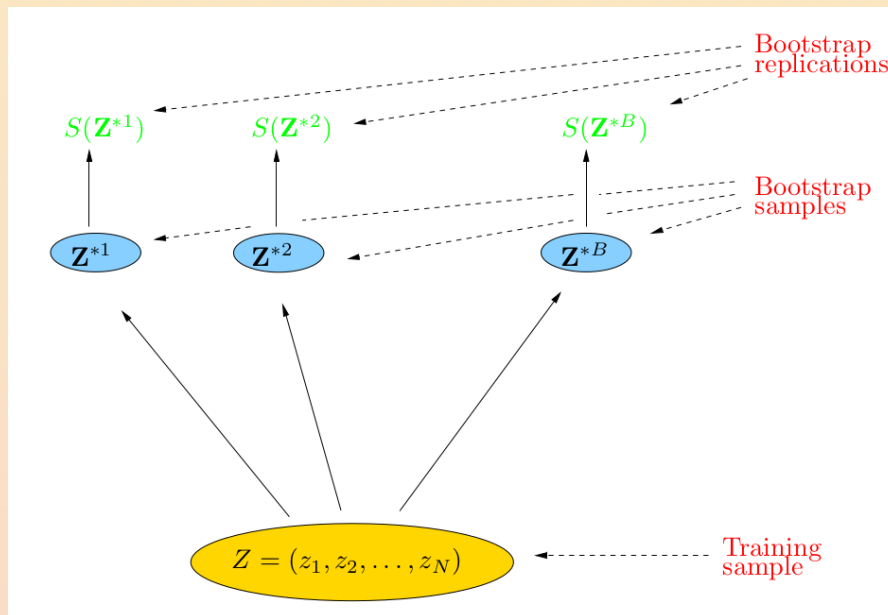
Model assessment: focus on prediction error

- **Cross-validation**: useful when data are scarce
 - Split the data into K parts ("folds")
 - For the k th part, fit the model to the other $K-1$ folds, and calculate test error as error on predicting the k th part data
 - Do this for all k , then combine the K estimates of the prediction error
 - Choose K
 - $K=N$ (leave-one-out), unbiased but high variance (training sets are basically the same)
 - Low K (5--10): Lower variance, but maybe bias (folds not representative of the data set)



Model assessment: focus on prediction error

- **Bootstrap**: a general tool to assessing statistical accuracy
 - Estimate the variance on the statistic $S(Z)$ (Z are the data)
 - Can be used as model assessment tool, or to improve an estimator
 - **Bagging** to combine weak learners (ensemble learning)



Exercises location

The public repository is at:

https://github.com/vischia/data_science_school_igfae2024

- The `README.md` contains instructions to run:
 - Locally on your machine (conda, virtualenv)
 - On Google Colab (in case you suspect your laptop may not be powerful enough for training neural networks in a timely manner (mostly relevant for tomorrow, today you should be fine))

Thank you!

