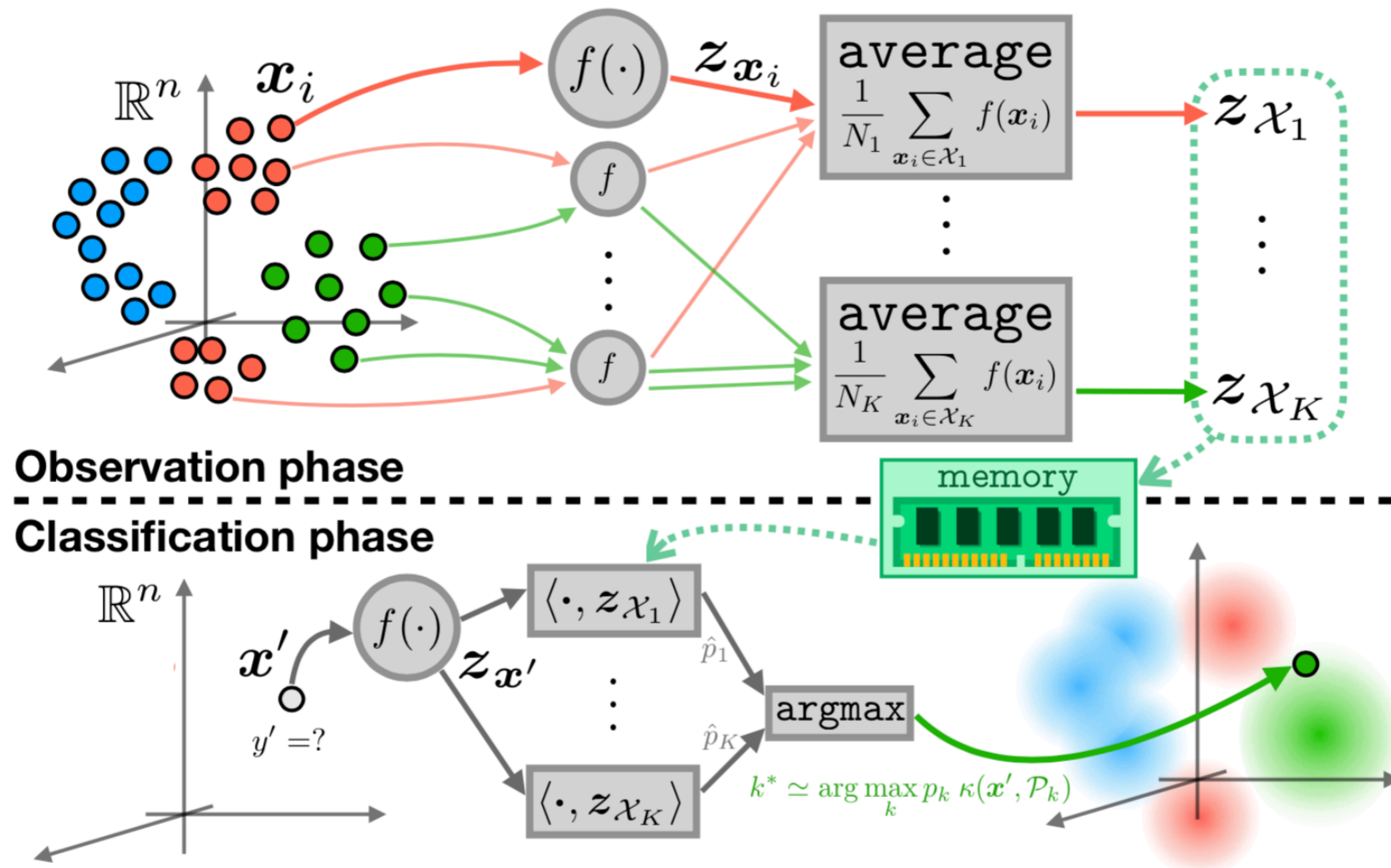


Compressive Classification

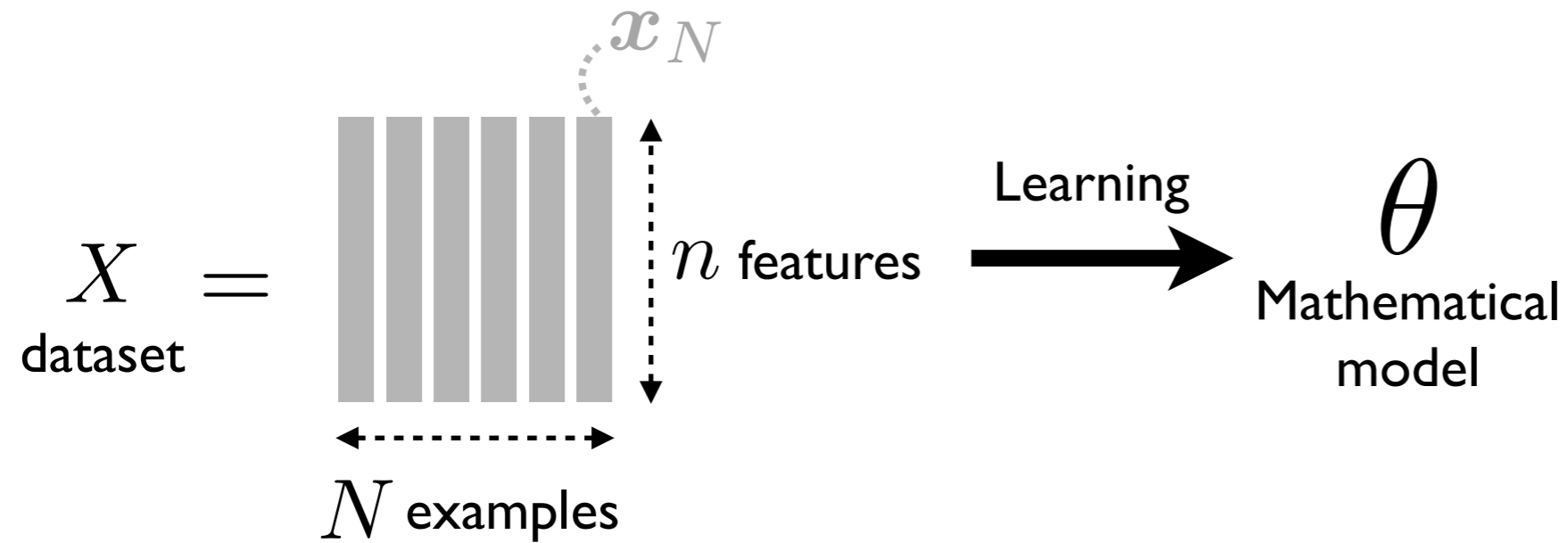
(Machine Learning without learning)



Vincent Schellekens

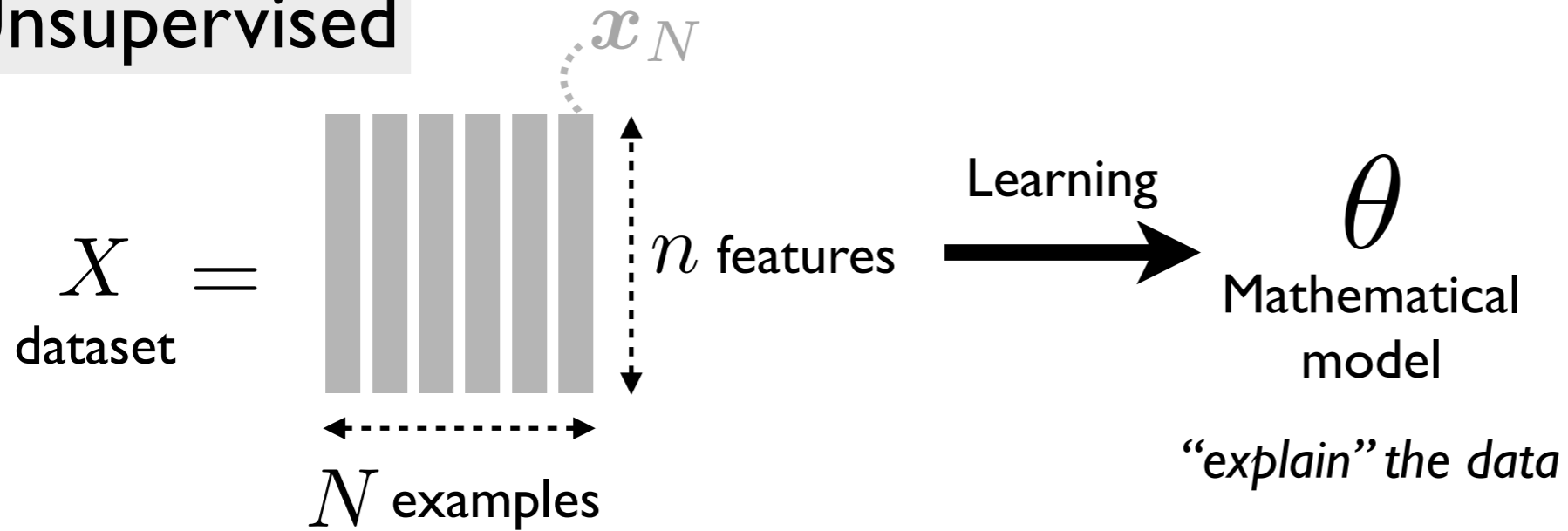
Laurent Jacques

Machine Learning



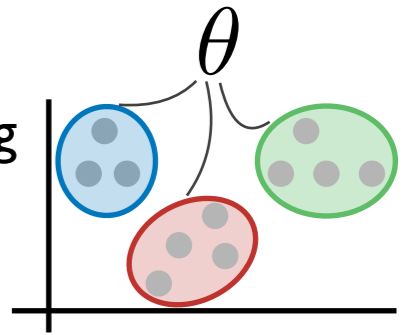
Machine Learning

Unsupervised

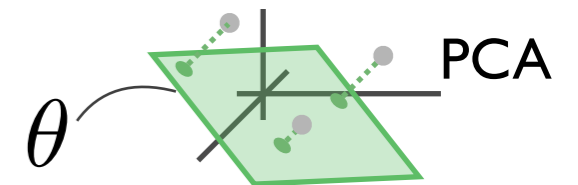


E.g.,

- Clustering



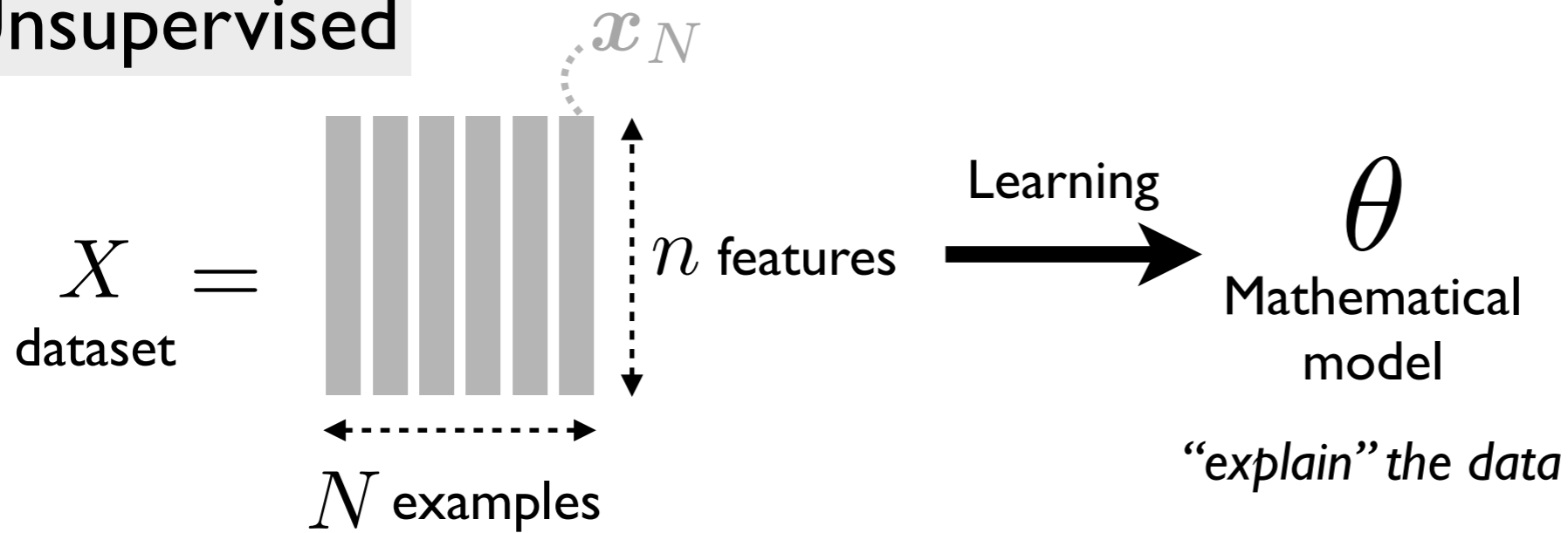
- Dimensionality reduction



- Autoencoder, GAN, SOM...

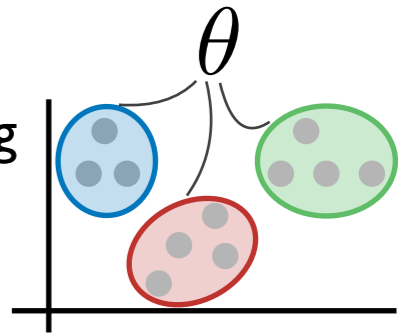
Machine Learning

Unsupervised

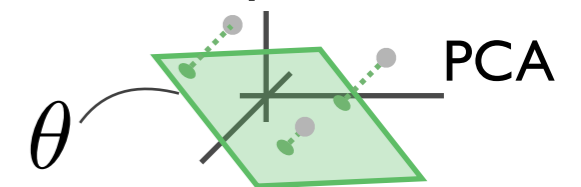


E.g.,

- Clustering

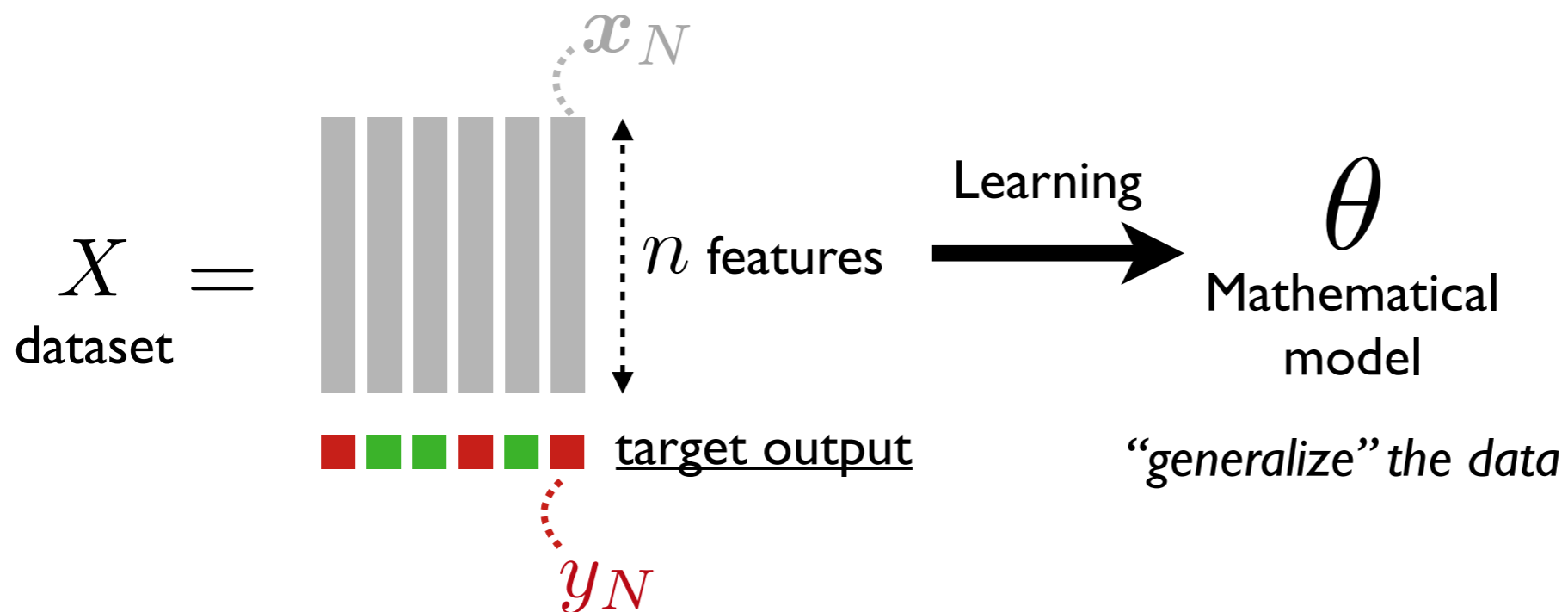


- Dimensionality reduction



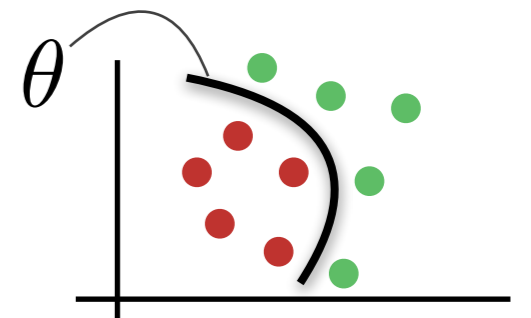
- Autoencoder, GAN, SOM...

Supervised

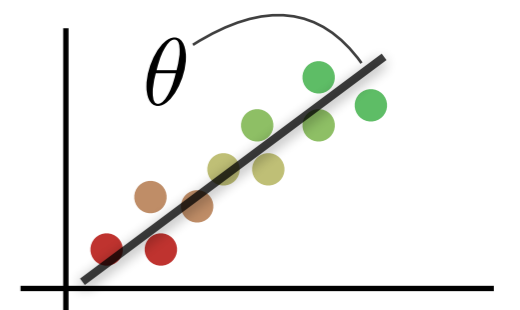


E.g.,

- Classification

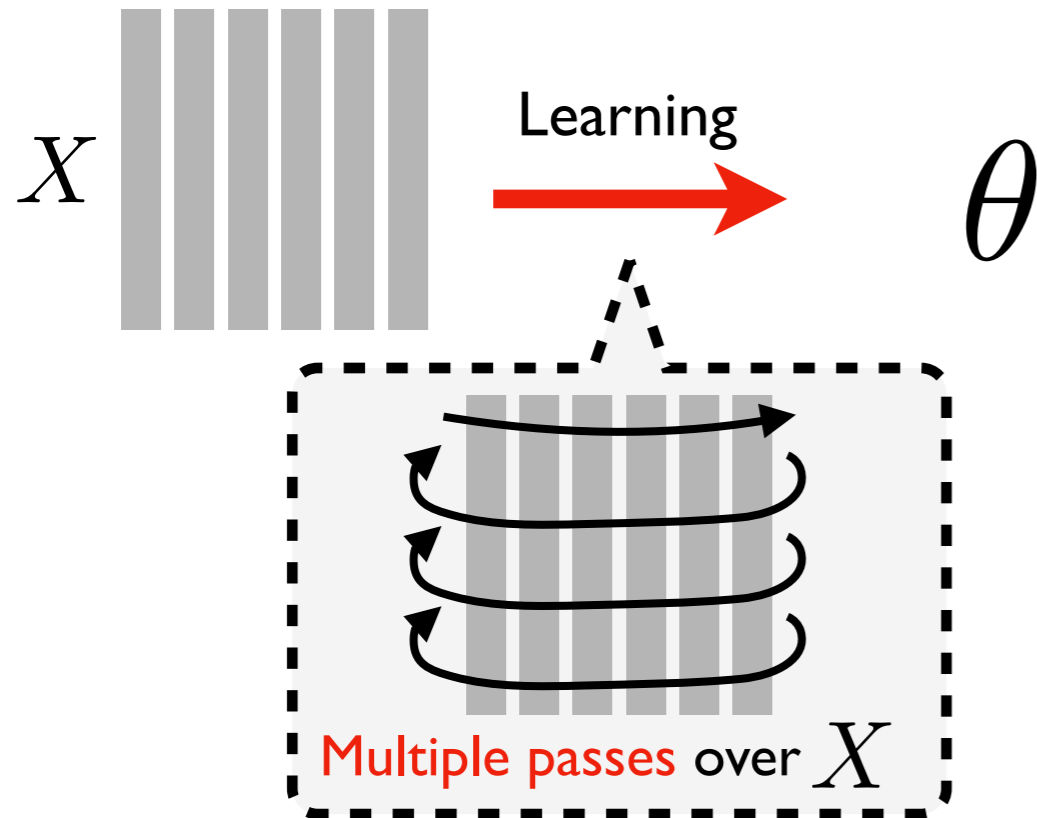


- Regression



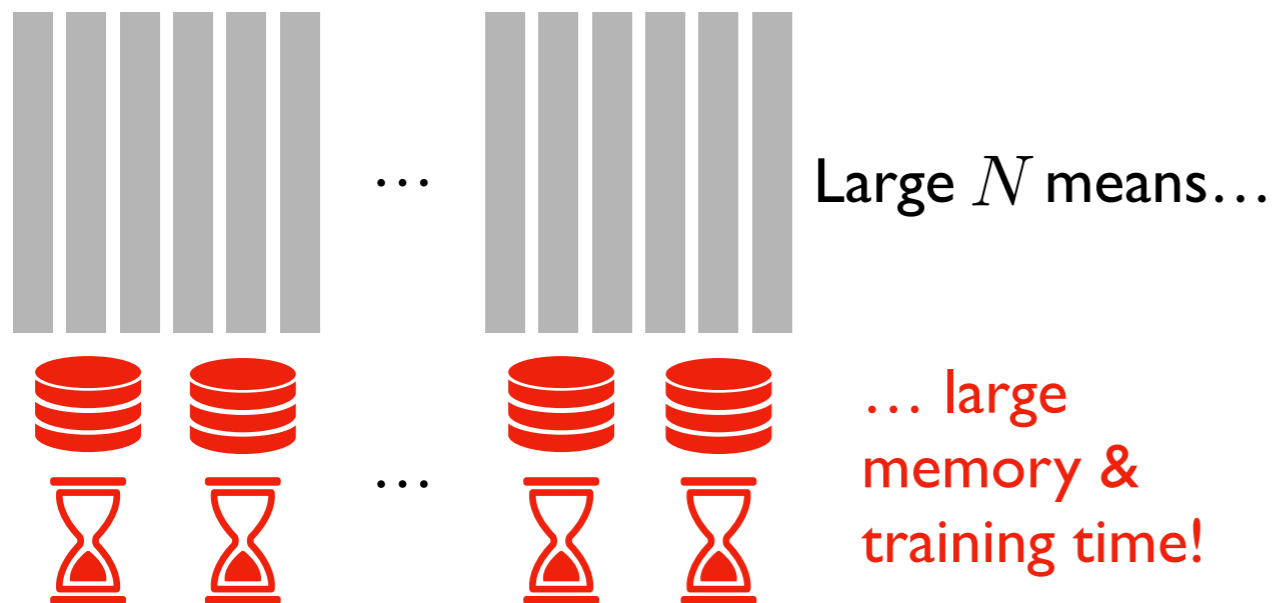
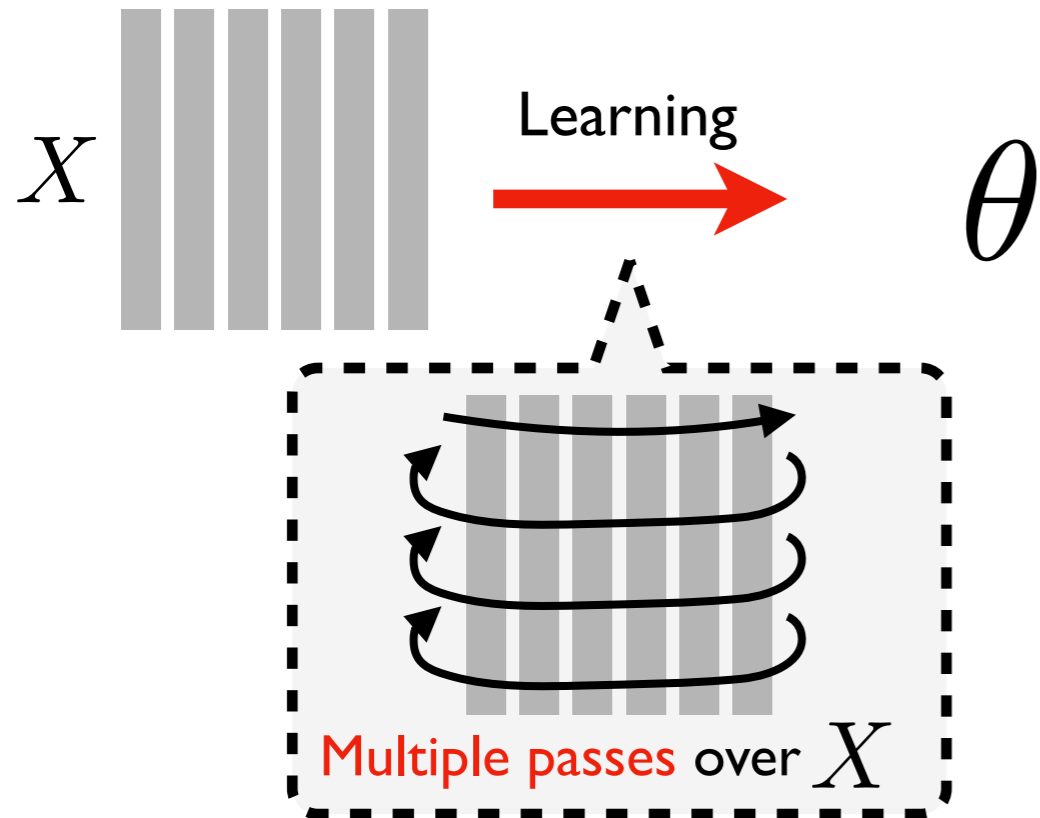
Compressive Learning

Usual machine learning



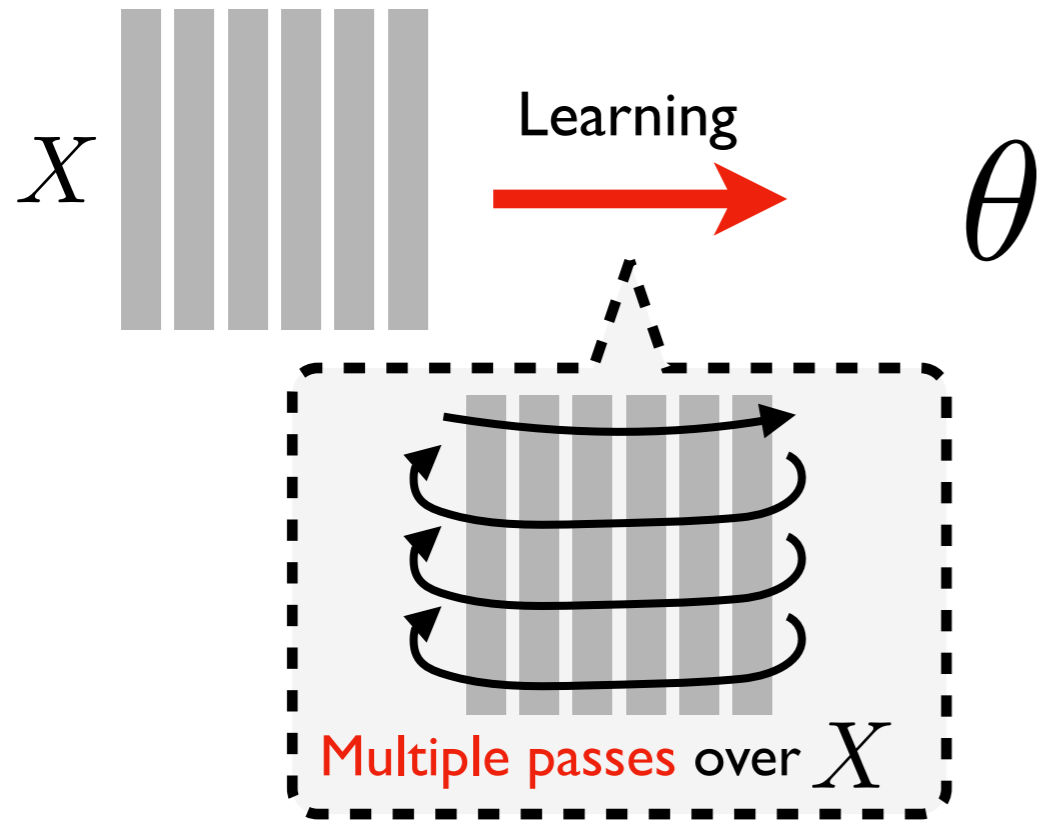
Compressive Learning

Usual machine learning

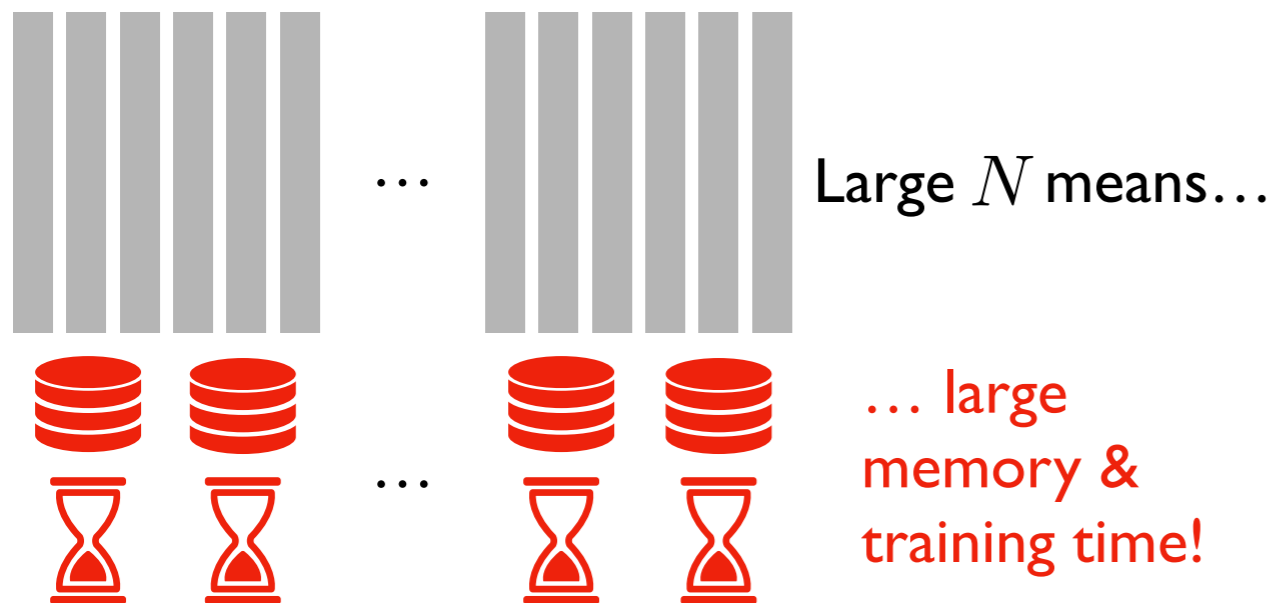
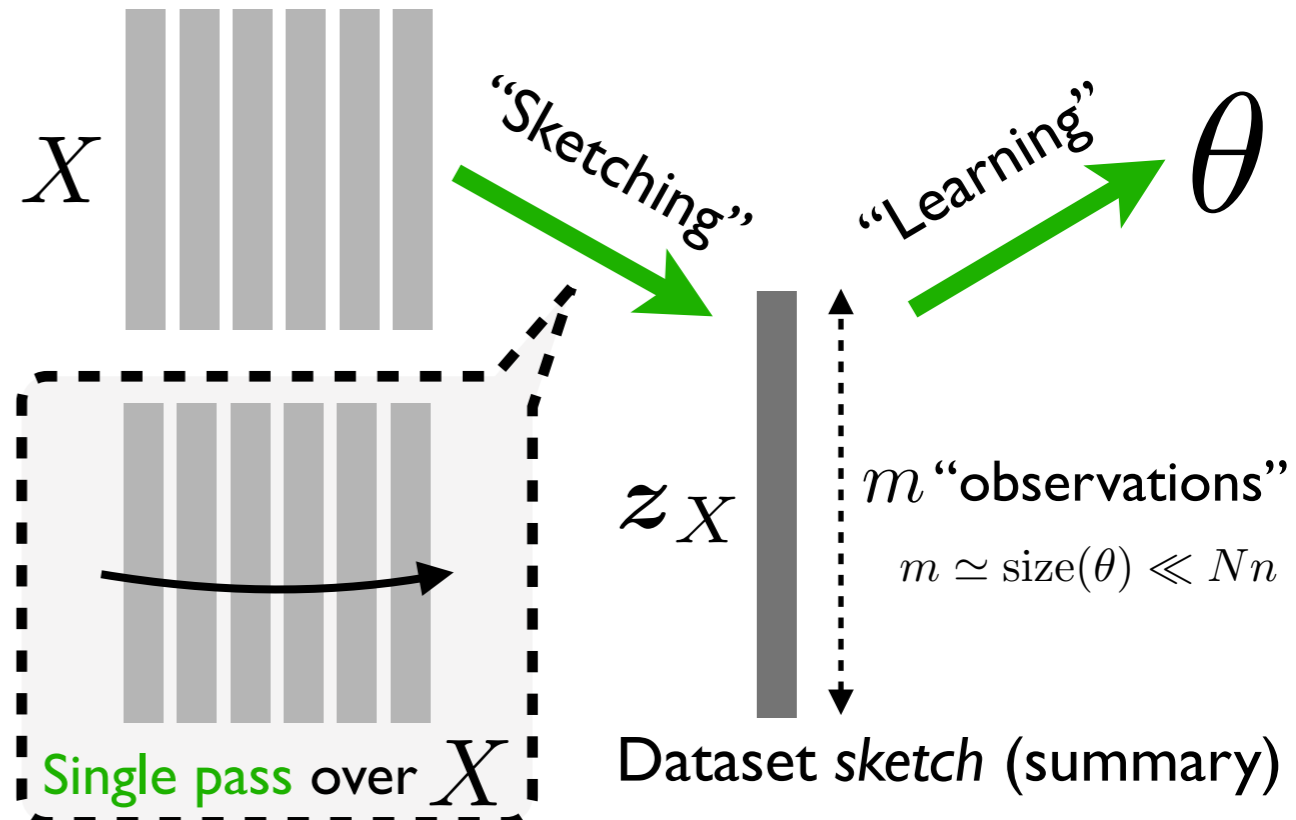


Compressive Learning

Usual machine learning

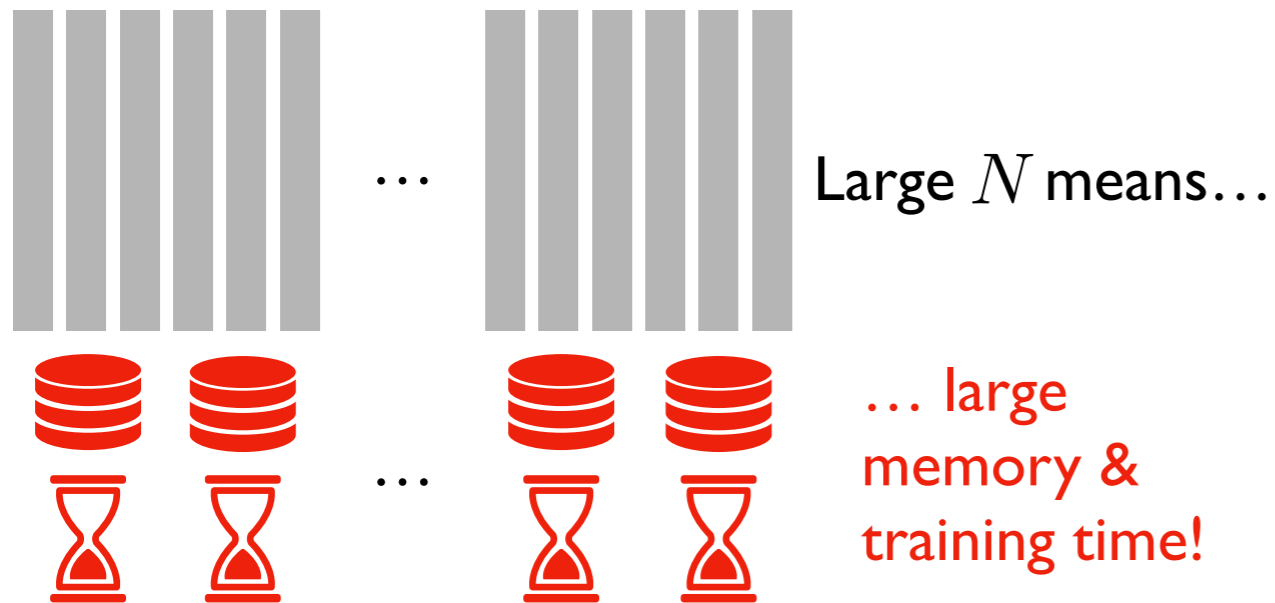
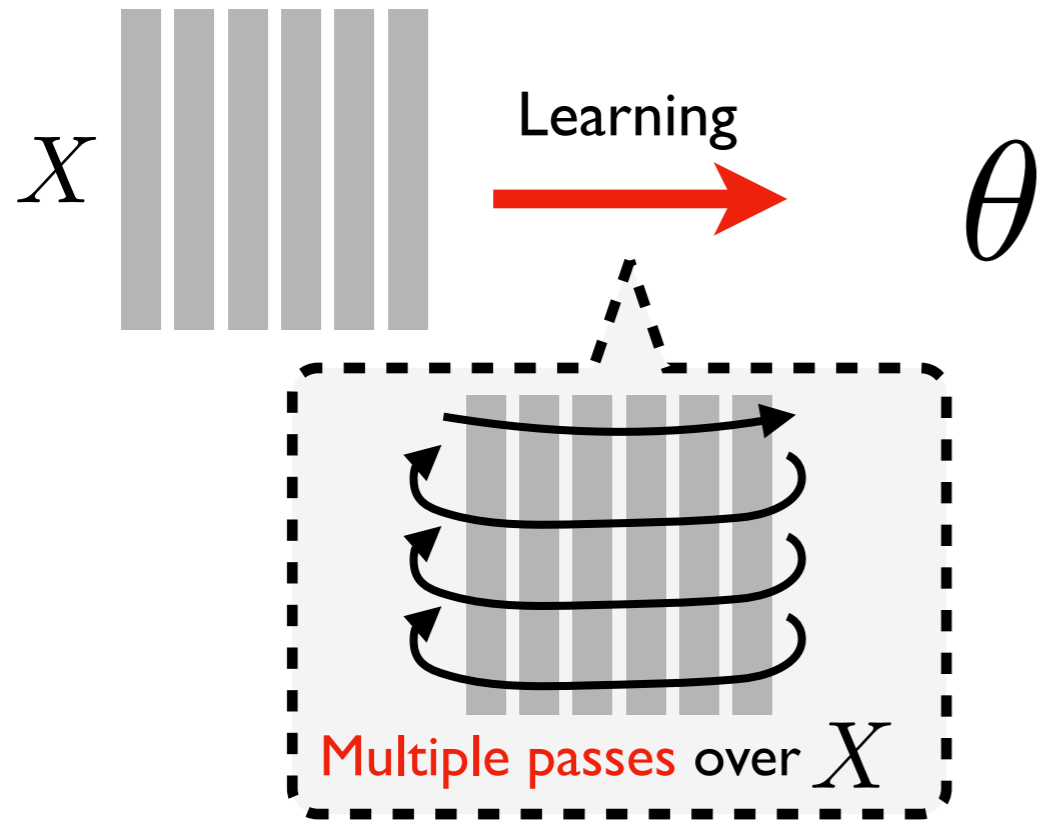


Compressive Learning e.g., [Gribonval-CL]

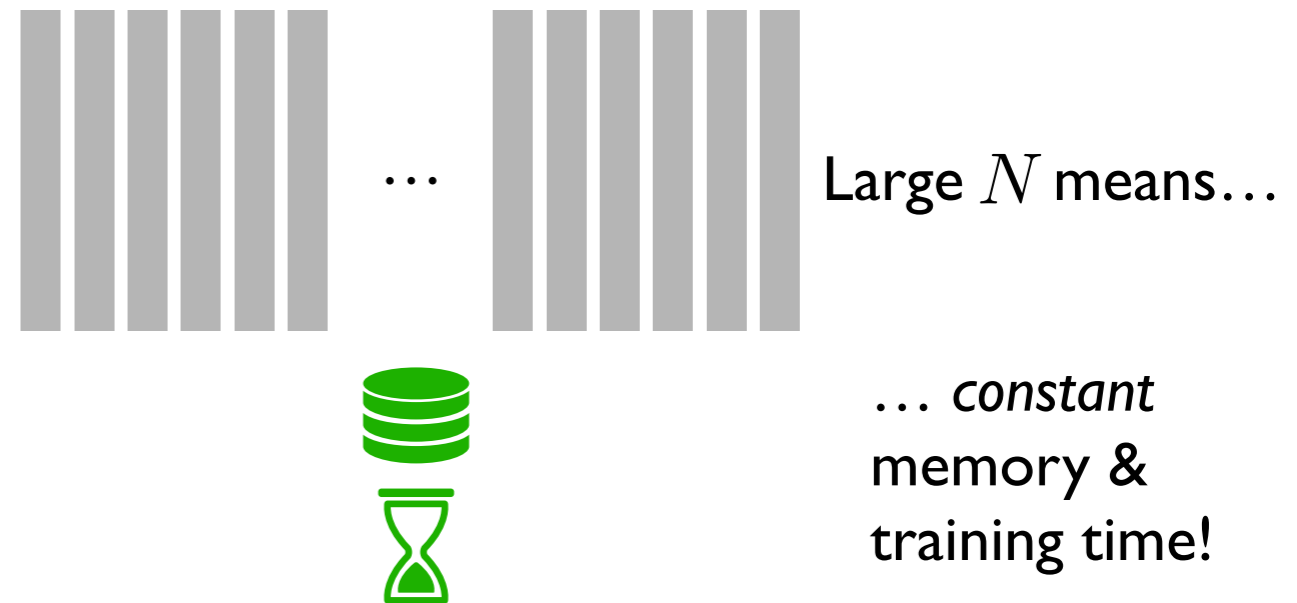
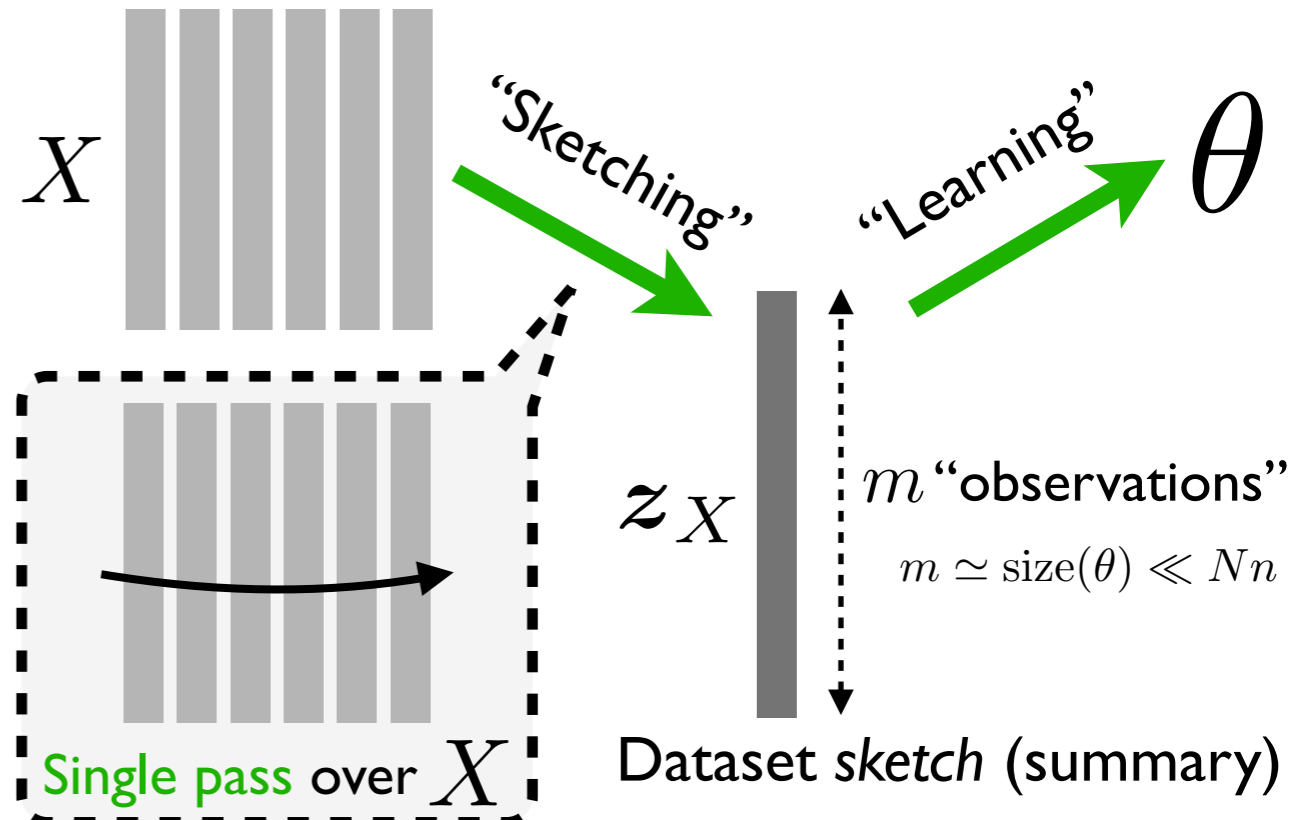


Compressive Learning

Usual machine learning

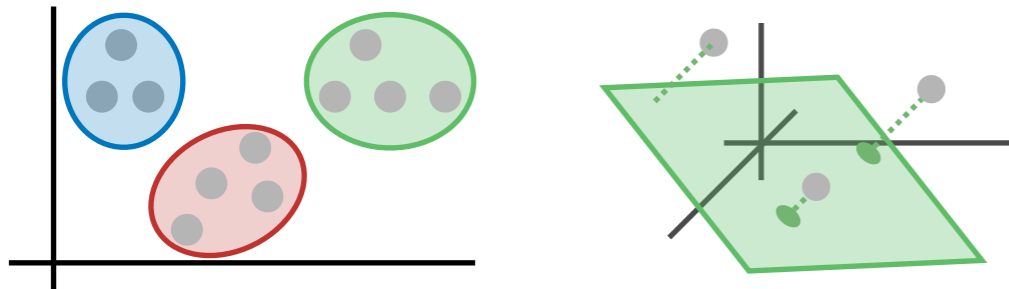


Compressive Learning e.g., [Gribonval-CL]



Previously on Compressive Learning...

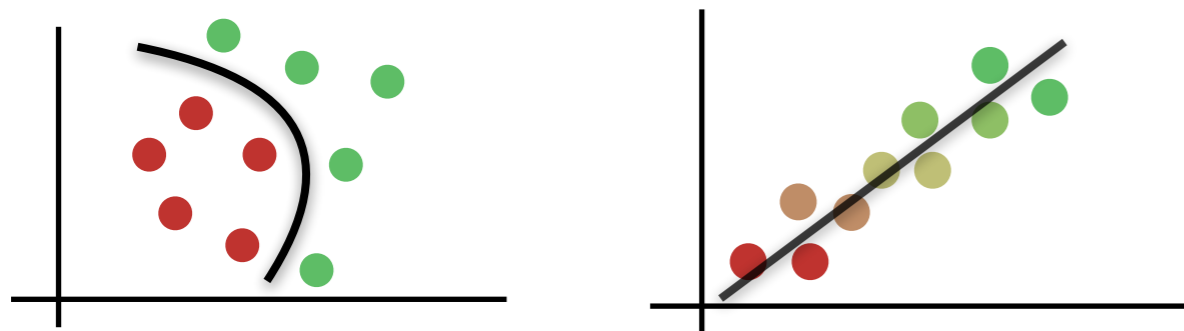
Unsupervised ML



Unsupervised Compressive Learning

- Compressive K-Means [Keriven-CKM]
- Compressive GMM estimation [Keriven-GMM]
- Compressive PCA [Gribonval-CL]

Supervised ML

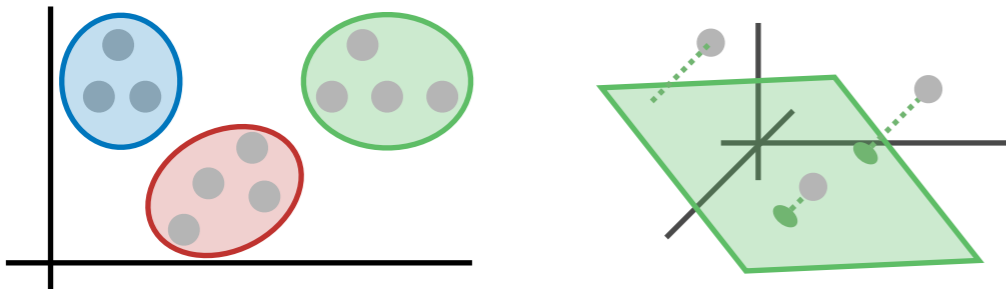


Supervised Compressive Learning

????

In this talk...

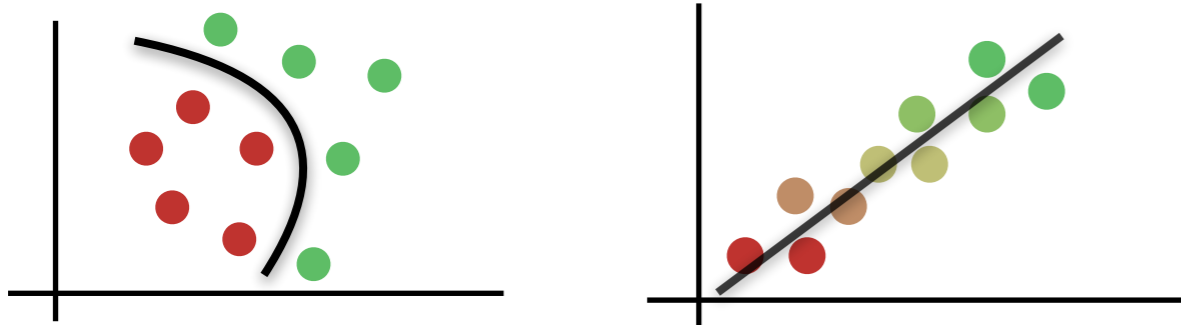
Unsupervised ML



Unsupervised Compressive Learning

- Compressive K-Means [Keriven-CKM]
- Compressive GMM estimation [Keriven-GMM]
- Compressive PCA [Gribonval-CL]

Supervised ML

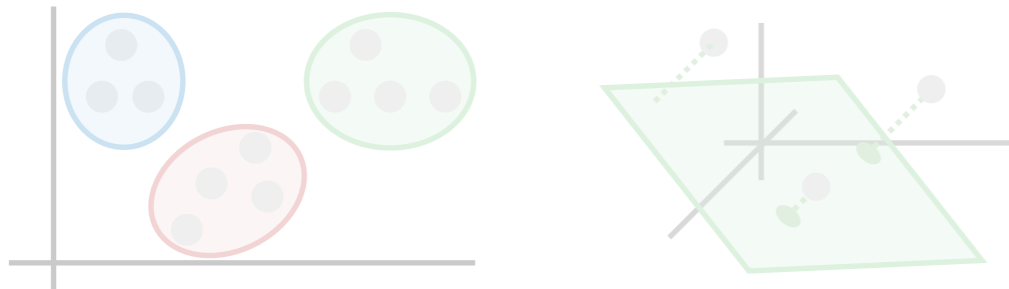


Supervised Compressive Learning

Compressive Classification
(a proof of concept)

In this talk...

Unsupervised ML

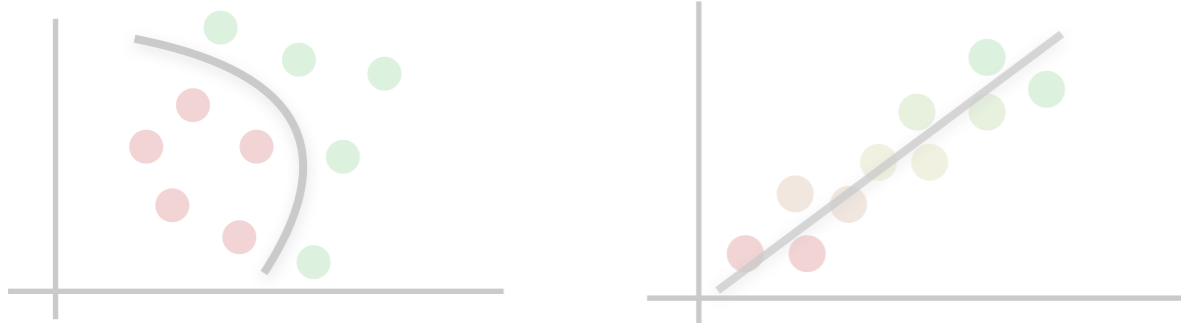


Unsupervised
Compressive Learning

- Compressive K-Means [Keriven-CKM]
- Compressive GMM estimation [Keriven-GMM]
- Compressive PCA [Gribonval-CL]

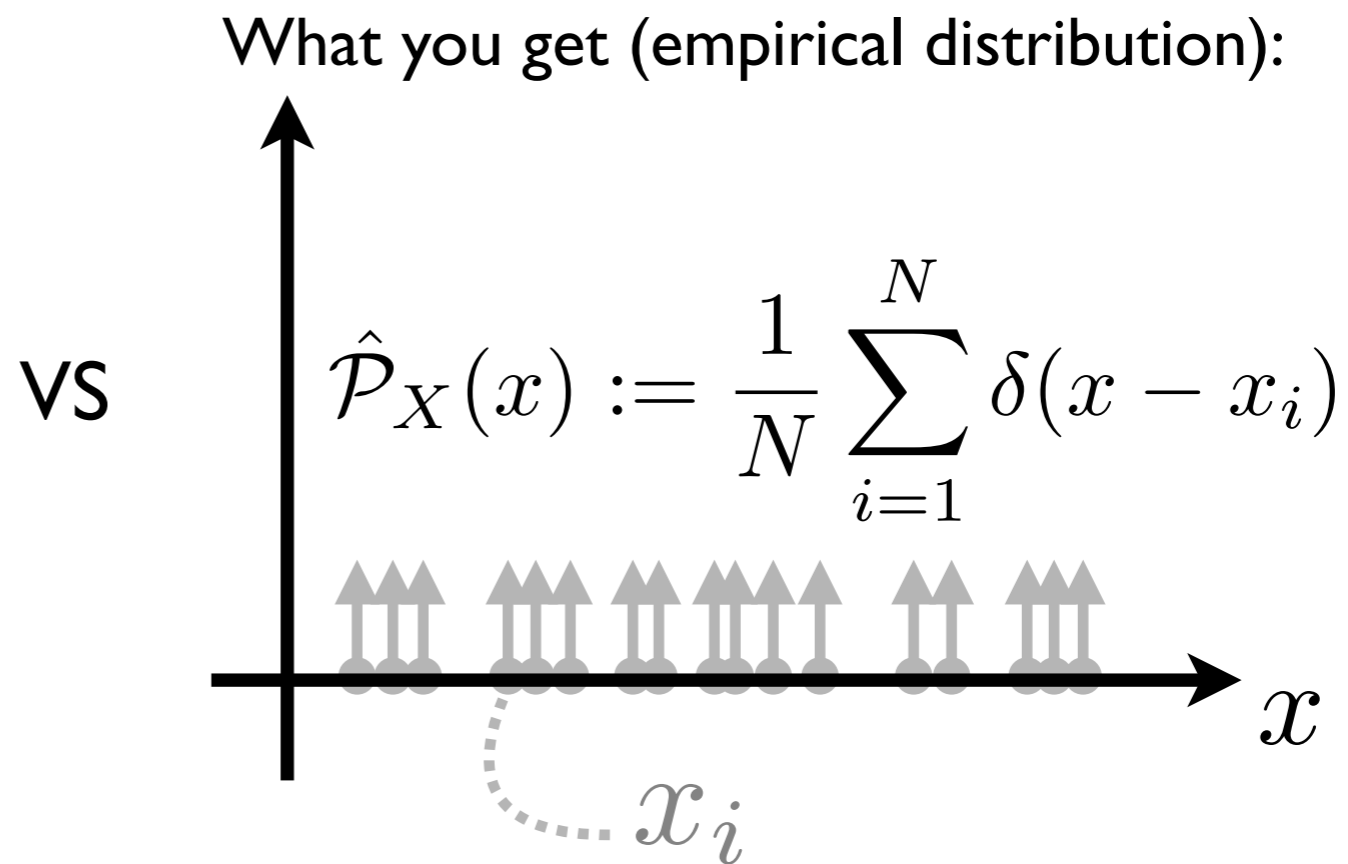
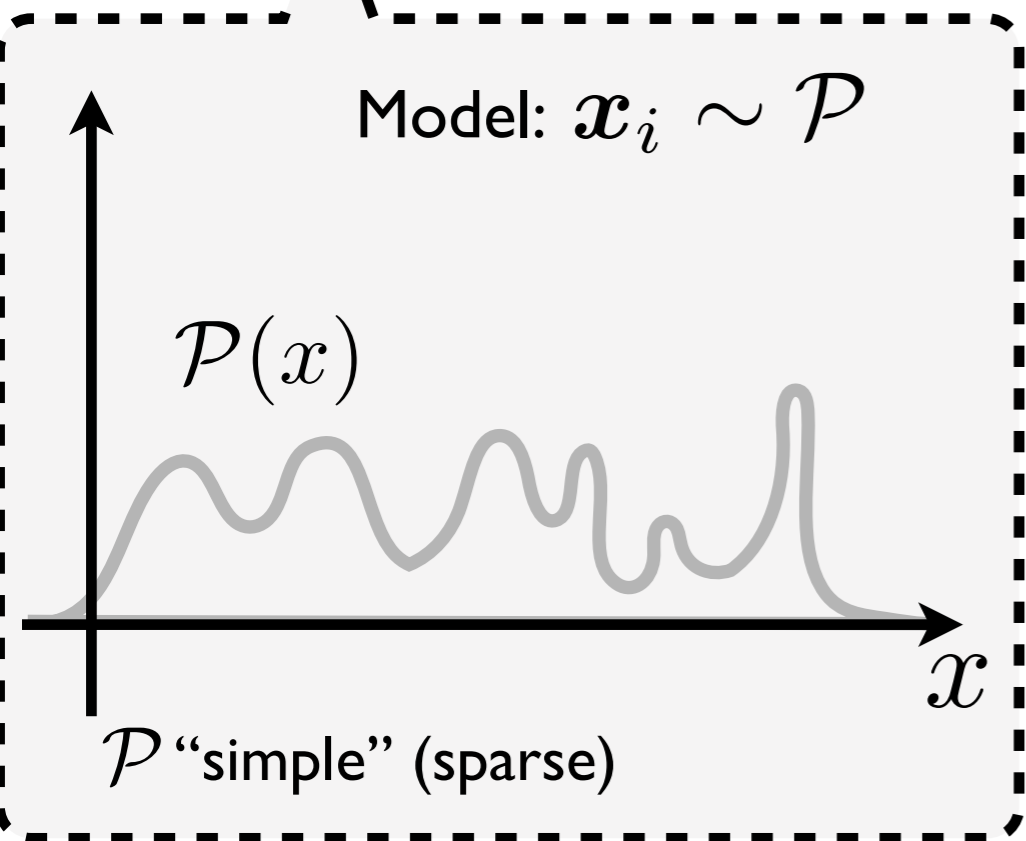
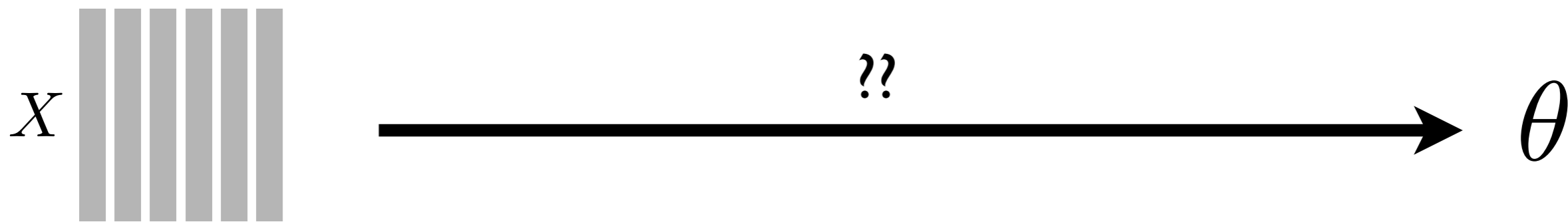
Preliminary 1: (Unsupervised) Compressive Learning Basics

Compressive Learning

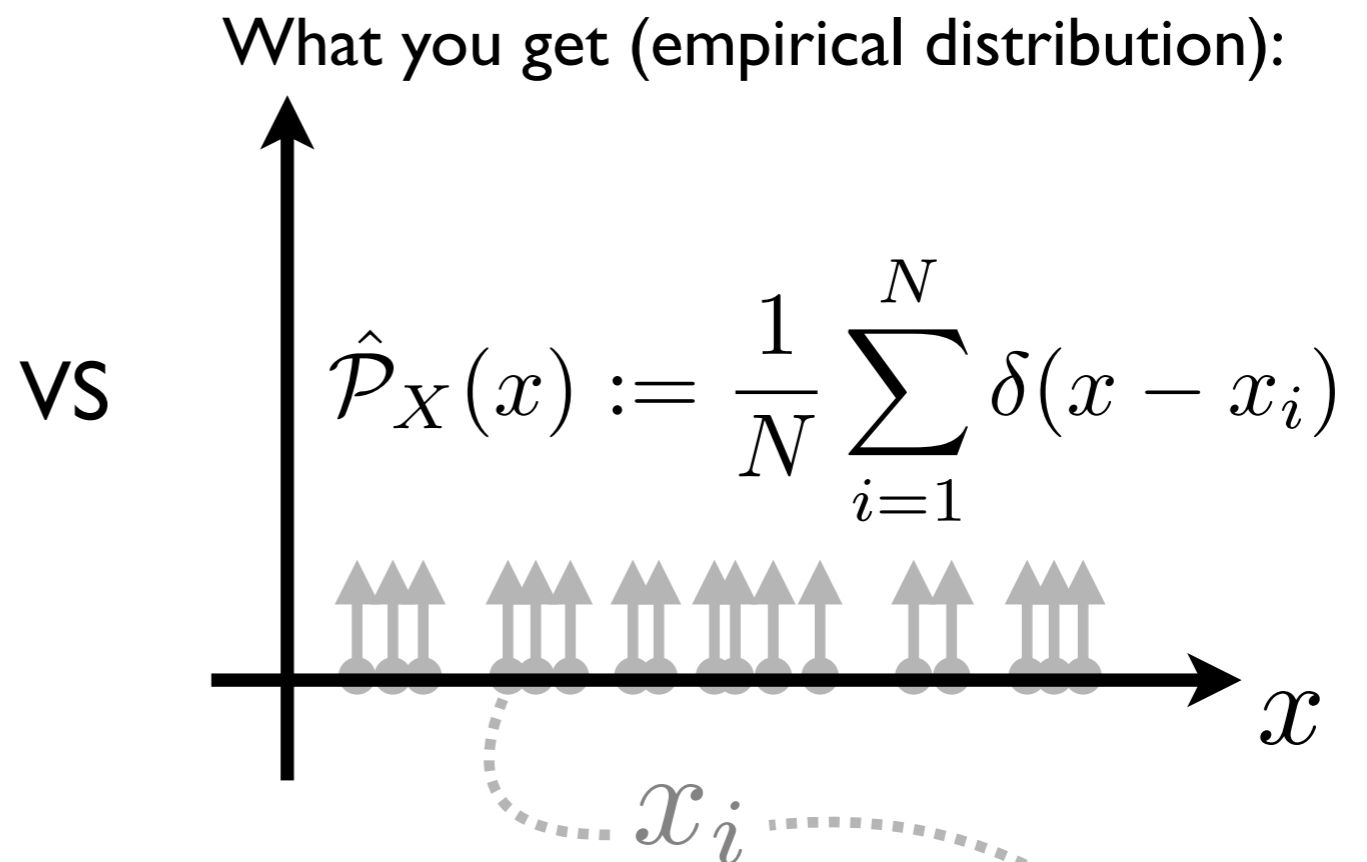
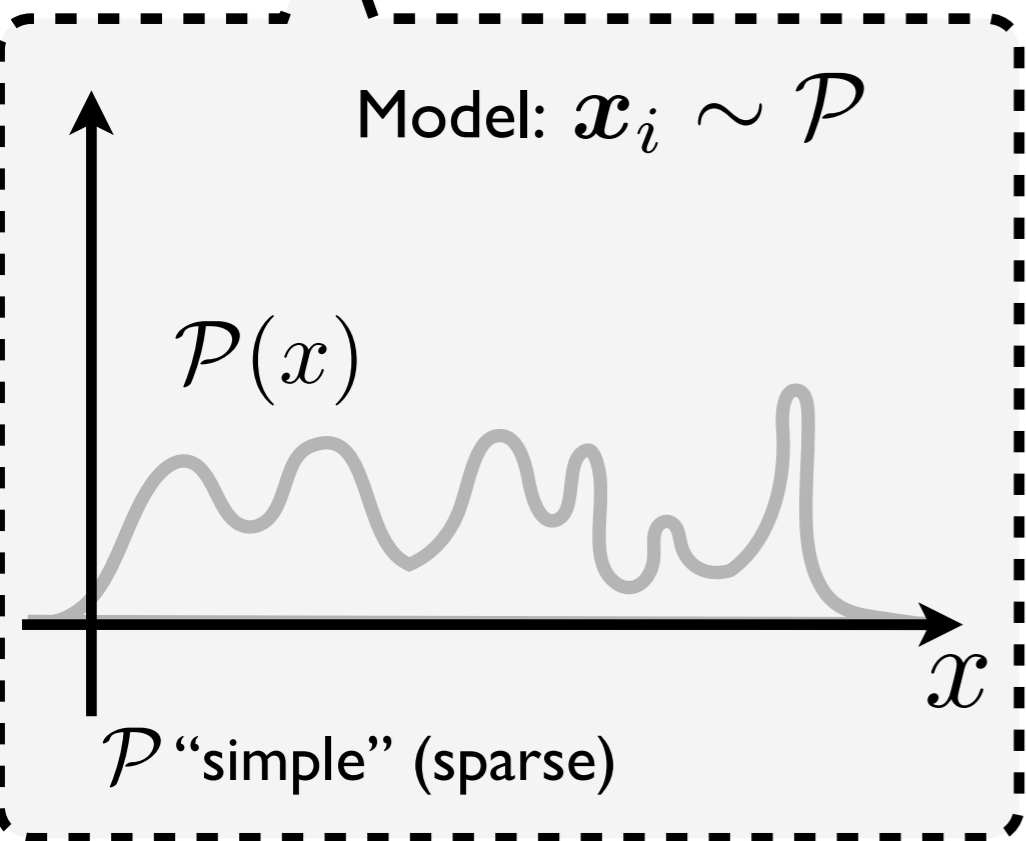
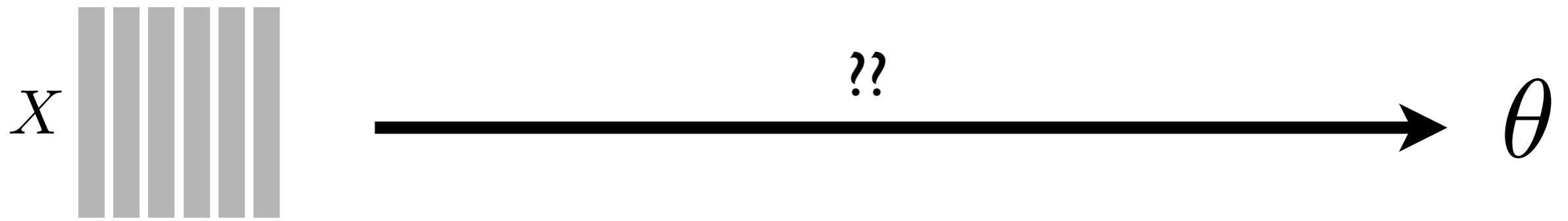


Compressive Classification
(a proof of concept)

Compressing a dataset

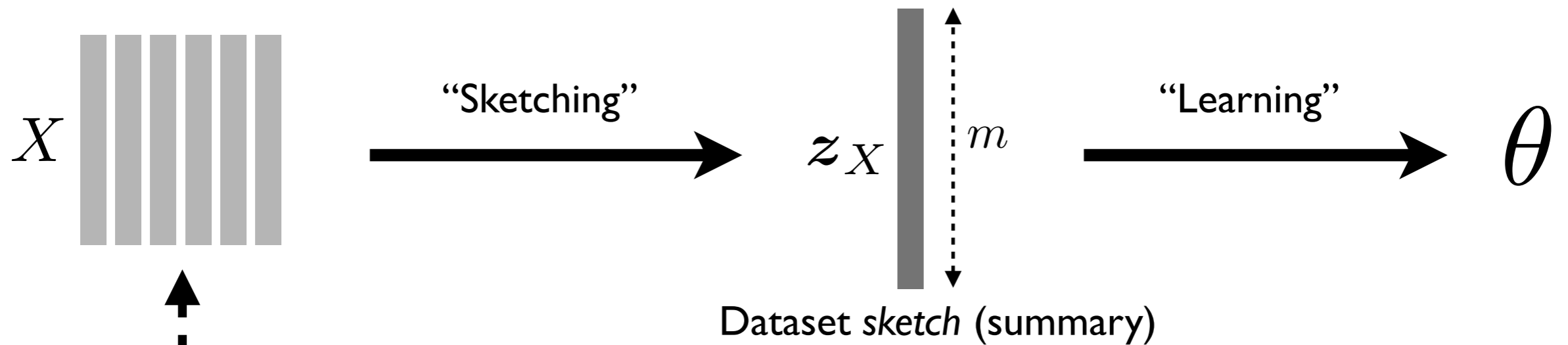


Compressing a **data distribution**

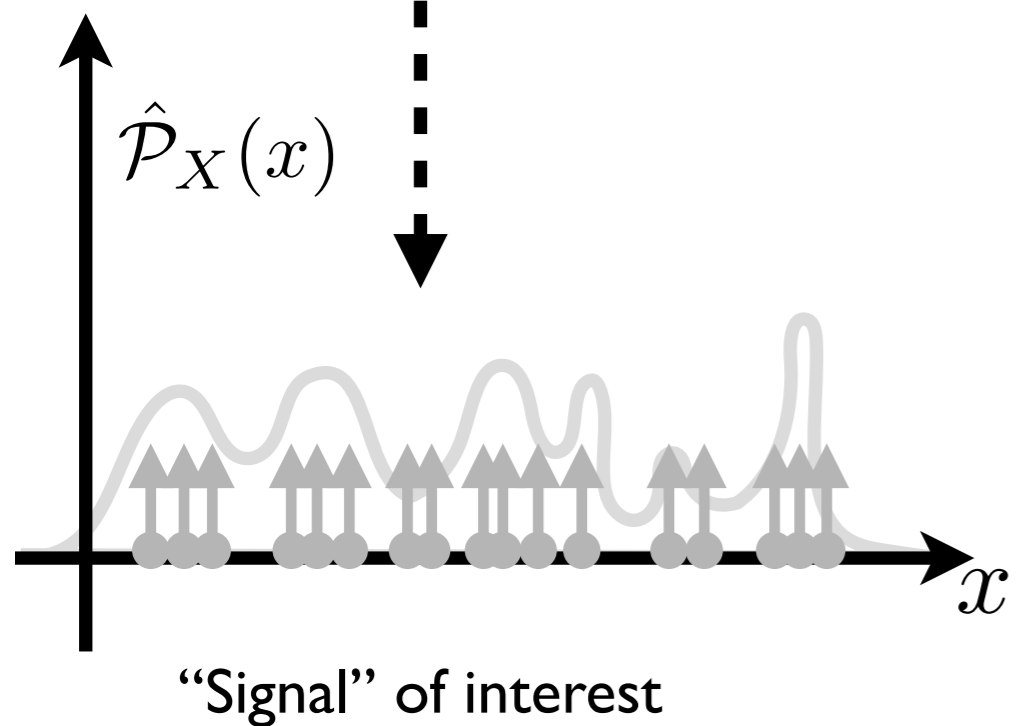


We don't care about individual realisations
Object of interest: the distribution \mathcal{P} !

Compressing a dataset distribution

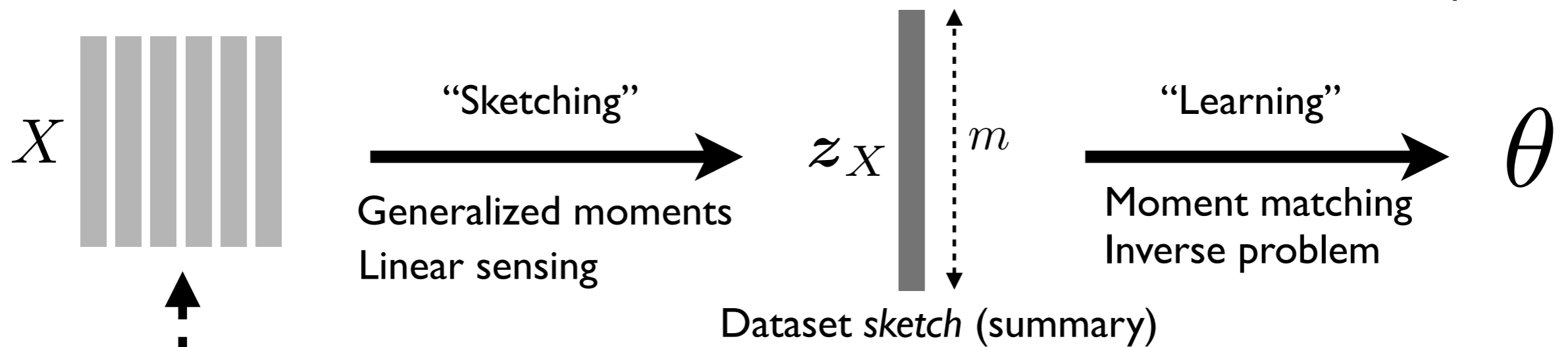


Dataset sketch (summary)



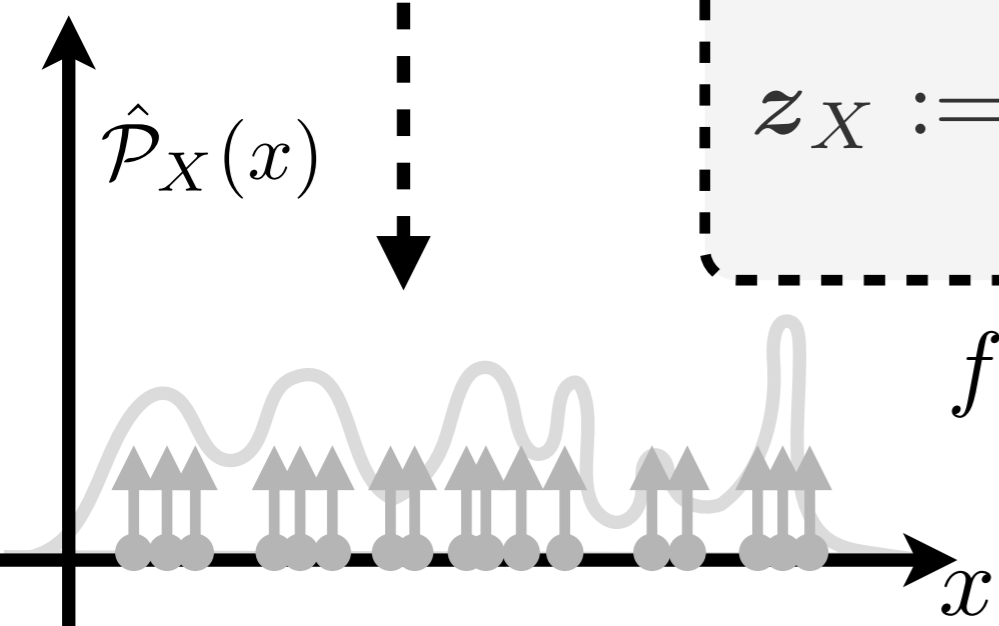
Compressing a **dataset** distribution

CS-inspired!



$$\begin{aligned}
 \mathcal{A}(\mathcal{P}) &:= \mathbb{E}_{\mathbf{x} \sim \mathcal{P}} f(\mathbf{x}) \\
 &\approx \\
 z_X &:= \mathcal{A}(\hat{\mathcal{P}}_X) = \frac{1}{N} \sum_{\mathbf{x}_i \in X} f(\mathbf{x}_i) \in \mathbb{C}^m
 \end{aligned}$$

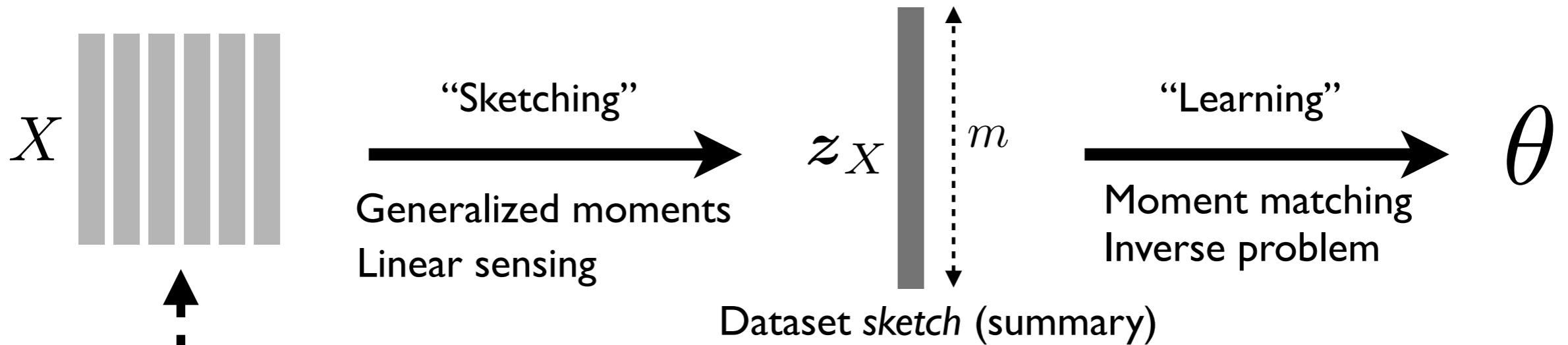
f computes m "features" of the data, to be averaged ("pooled")



"Signal" of interest

Compressing a **dataset** distribution

CS-inspired!

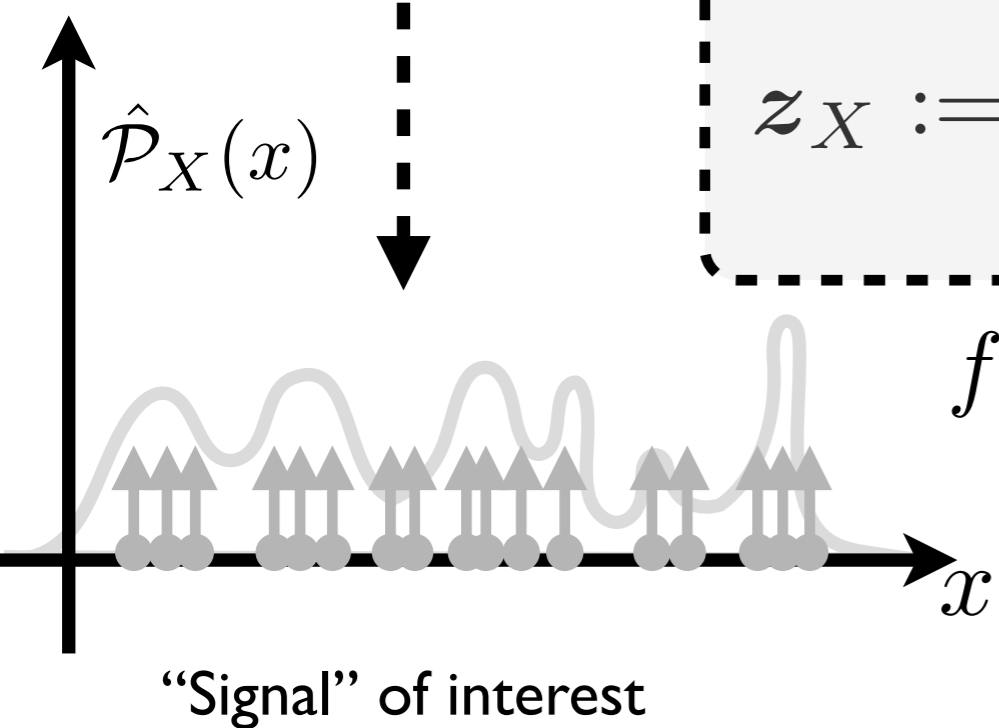


$$\begin{aligned}
 \mathcal{A}(\mathcal{P}) &:= \mathbb{E}_{\mathbf{x} \sim \mathcal{P}} f(\mathbf{x}) \\
 &\approx z_X := \mathcal{A}(\hat{\mathcal{P}}_X) = \frac{1}{N} \sum_{\mathbf{x}_i \in X} f(\mathbf{x}_i) \in \mathbb{C}^m
 \end{aligned}$$

f computes m "features" of the data, to be averaged ("pooled")

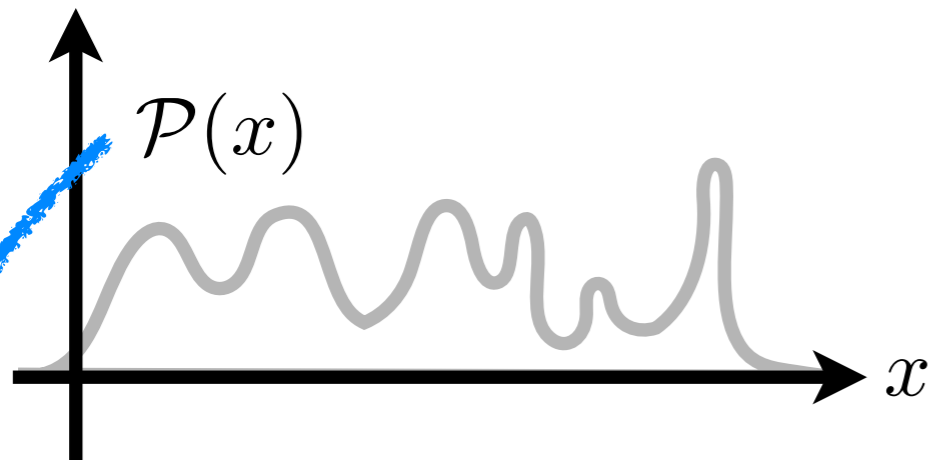
E.g. Random Fourier Features

$$f(\mathbf{x})_{\text{RFF}} = \left[\exp(i\omega_j^T \mathbf{x}) \right]_{j=1}^m$$



Geometric interpretation

[Smola]

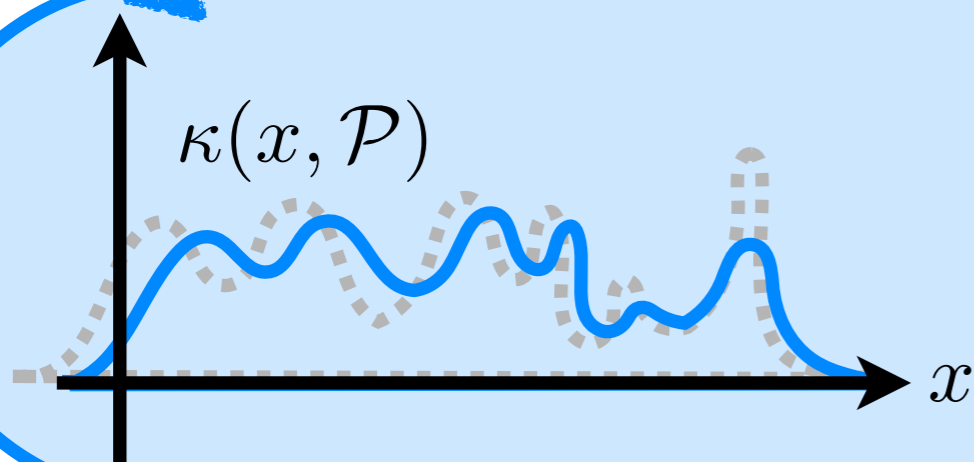


“Mean Map”

$$\kappa(\cdot, \mathcal{P}) := \mathbb{E}_{\mathbf{x}' \sim \mathcal{P}} \kappa(\cdot, \mathbf{x}')$$

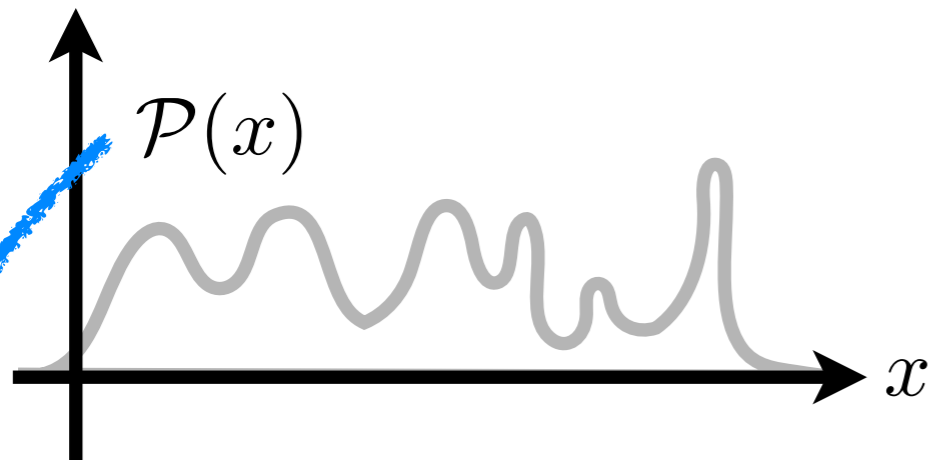
f approximates a kernel
 $\langle f(\mathbf{x}), f(\mathbf{x}') \rangle \simeq \kappa(\mathbf{x}, \mathbf{x}')$
associated with a RKHS \mathcal{H}_κ

“ $\mathcal{H}_\kappa = \text{span}(\{\kappa(\cdot, \mathbf{u})\})$ ”



\mathcal{H}_κ

Geometric interpretation

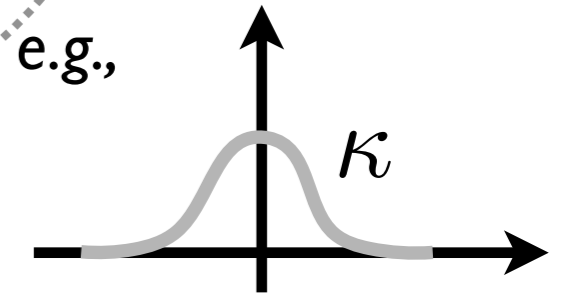


$$f(\mathbf{x})_{\text{RFF}} = \left[\exp(i\omega_j^T \mathbf{x}) \right]_{j=1}^m \quad \text{with } \omega_j \sim \Lambda :$$

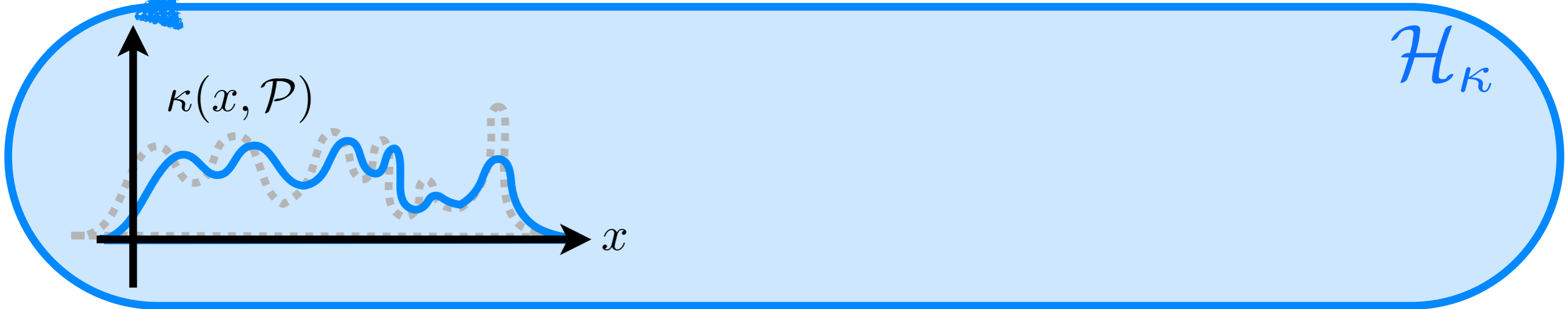
$$\langle f(\mathbf{x})_{\text{RFF}}, f(\mathbf{x}')_{\text{RFF}} \rangle \simeq \kappa(\mathbf{x}, \mathbf{x}') = K(\mathbf{u} = \mathbf{x} - \mathbf{x}') = (F\Lambda)(\mathbf{u})$$

“Mean Map”
 $\kappa(\cdot, \mathcal{P}) := \mathbb{E}_{\mathbf{x}' \sim \mathcal{P}} \kappa(\cdot, \mathbf{x}')$

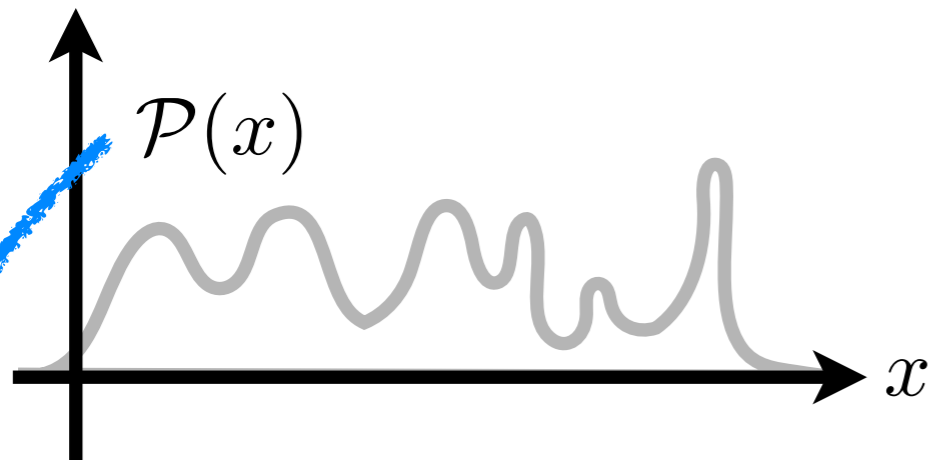
f approximates a kernel
 $\langle f(\mathbf{x}), f(\mathbf{x}') \rangle \simeq \kappa(\mathbf{x}, \mathbf{x}')$
 associated with a RKHS \mathcal{H}_κ



“ $\mathcal{H}_\kappa = \text{span}(\{\kappa(\cdot, \mathbf{u})\})$ ”



Geometric interpretation



“Mean Map”

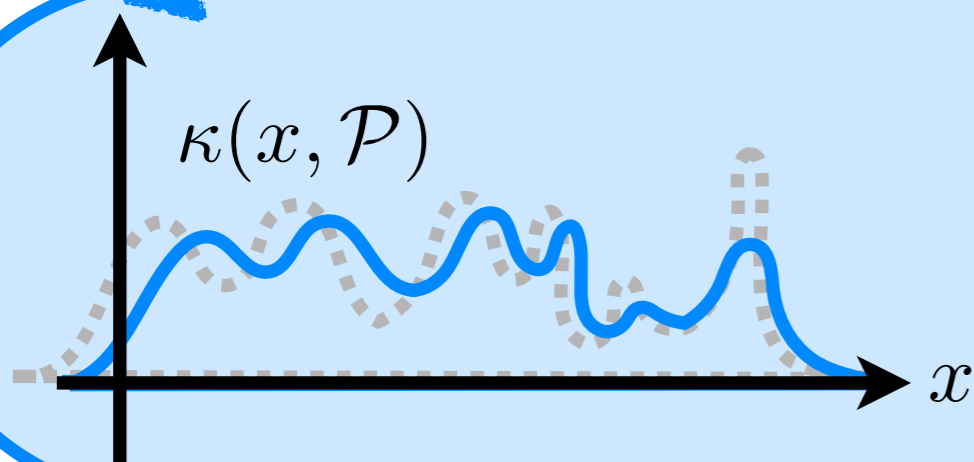
$$\kappa(\cdot, \mathcal{P}) := \mathbb{E}_{x' \sim \mathcal{P}} \kappa(\cdot, x')$$

f approximates a kernel
 $\langle f(\mathbf{x}), f(\mathbf{x}') \rangle \simeq \kappa(\mathbf{x}, \mathbf{x}')$
associated with a RKHS \mathcal{H}_κ

$m \rightarrow \infty$

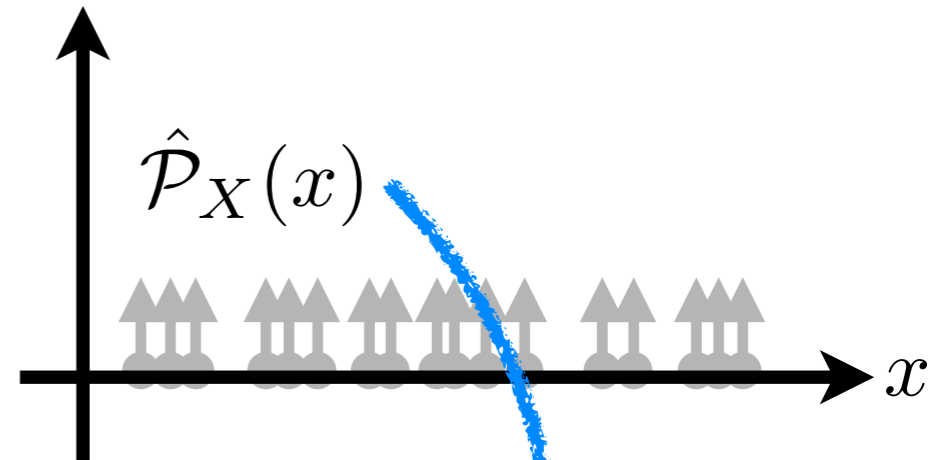
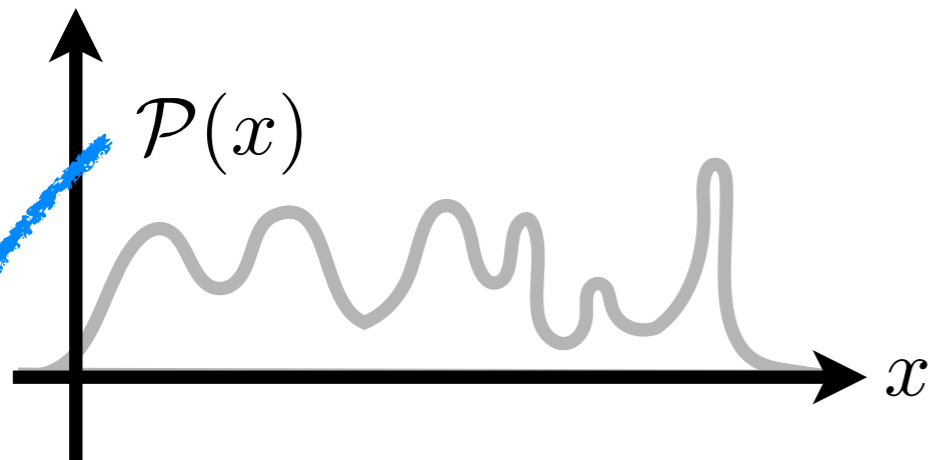
Key fact: the map $\mathcal{A} : \mathcal{P} \rightarrow \mathbb{E}_{x \sim \mathcal{P}} f(\mathbf{x}) \in \mathbb{C}^m$
approximatively preserves the geometry of \mathcal{H}_κ

“ $\mathcal{H}_\kappa = \text{span}(\{\kappa(\cdot, u)\})$ ”



\mathcal{H}_κ

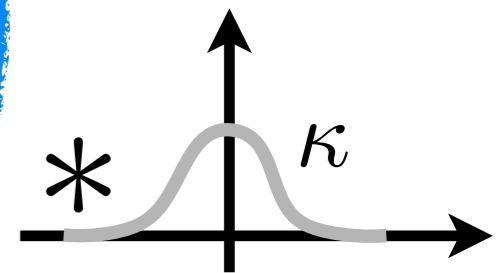
Geometric interpretation



“Mean Map”

$$\kappa(\cdot, \mathcal{P}) := \mathbb{E}_{x' \sim \mathcal{P}} \kappa(\cdot, x')$$

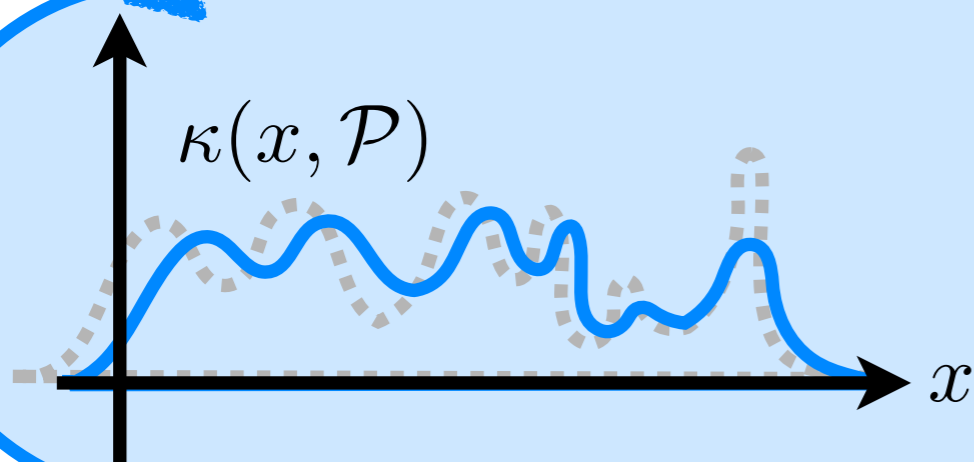
f approximates a kernel
 $\langle f(\mathbf{x}), f(\mathbf{x}') \rangle \simeq \kappa(\mathbf{x}, \mathbf{x}')$
 associated with a RKHS \mathcal{H}_κ



$m \rightarrow \infty$

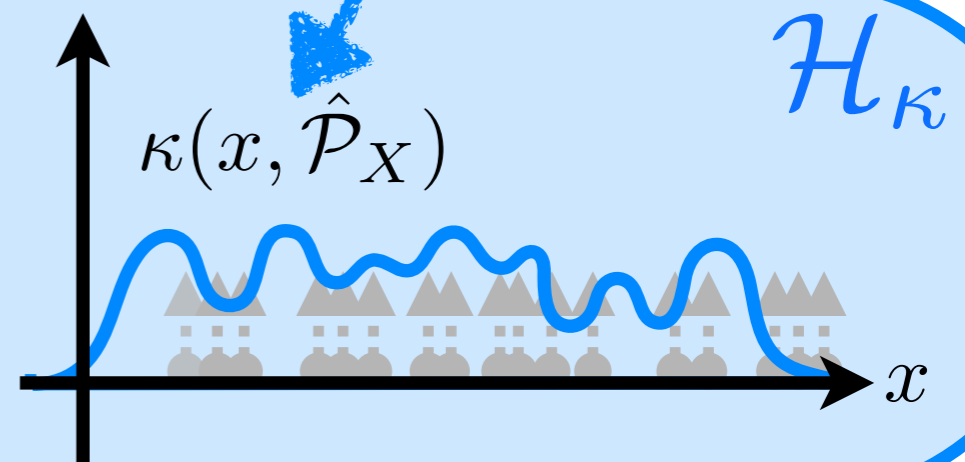
Key fact: the map $\mathcal{A} : \mathcal{P} \rightarrow \mathbb{E}_{x \sim \mathcal{P}} f(x) \in \mathbb{C}^m$
 approximatively preserves the geometry of \mathcal{H}_κ

“ $\mathcal{H}_\kappa = \text{span}(\{\kappa(\cdot, u)\})$ ”



$N \rightarrow \infty$

\simeq



In this talk...

Unsupervised ML

Unsupervised
Compressive Learning

- Compressive K-Means [Keriven-CKM]
- Compressive GMM [Keriven-GMM]
- Compressive CL [Keriven-CL]

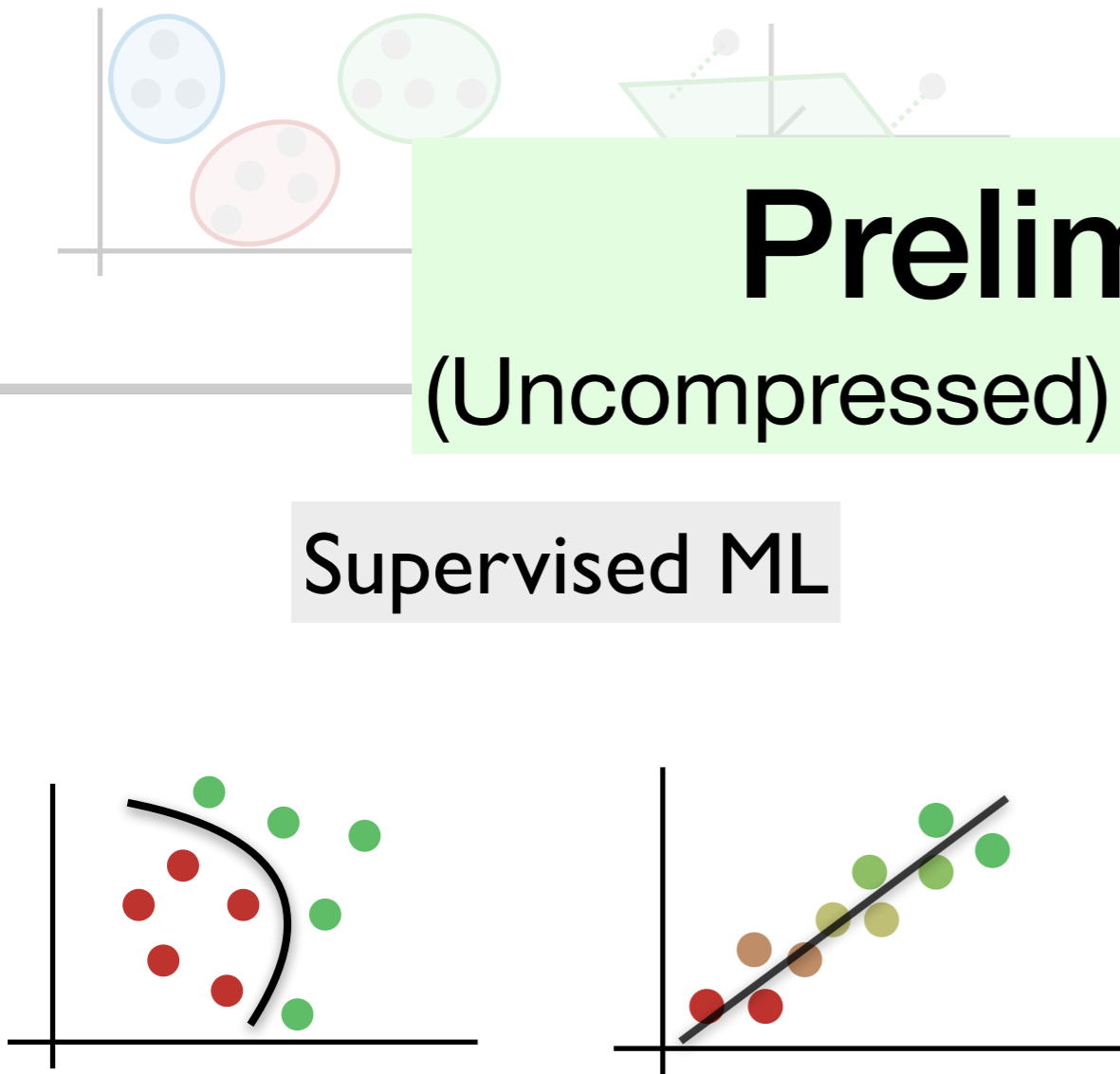
Preliminary 2:

(Uncompressed) Classification Basics

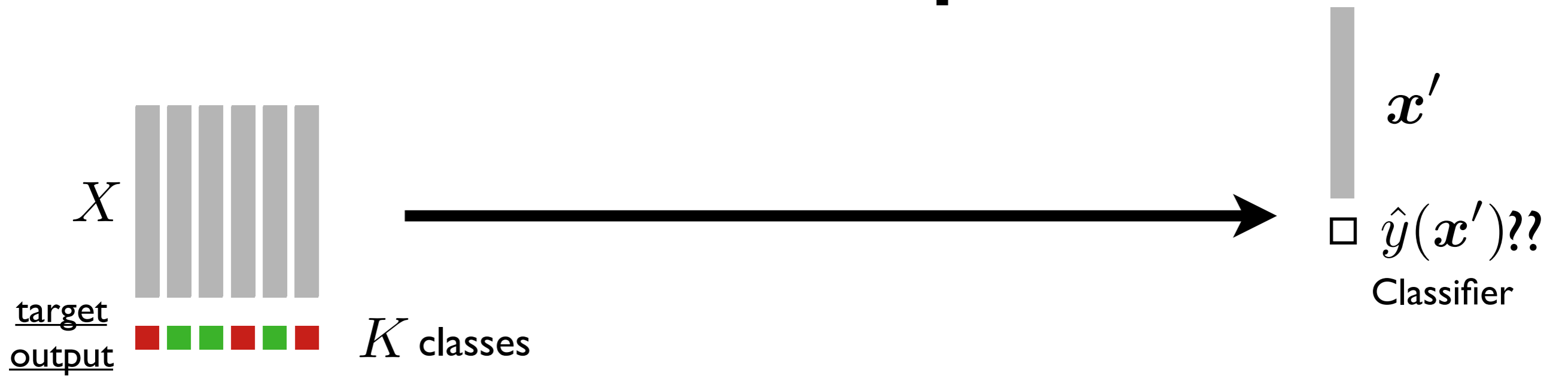
Supervised ML

Supervised
Compressive Learning

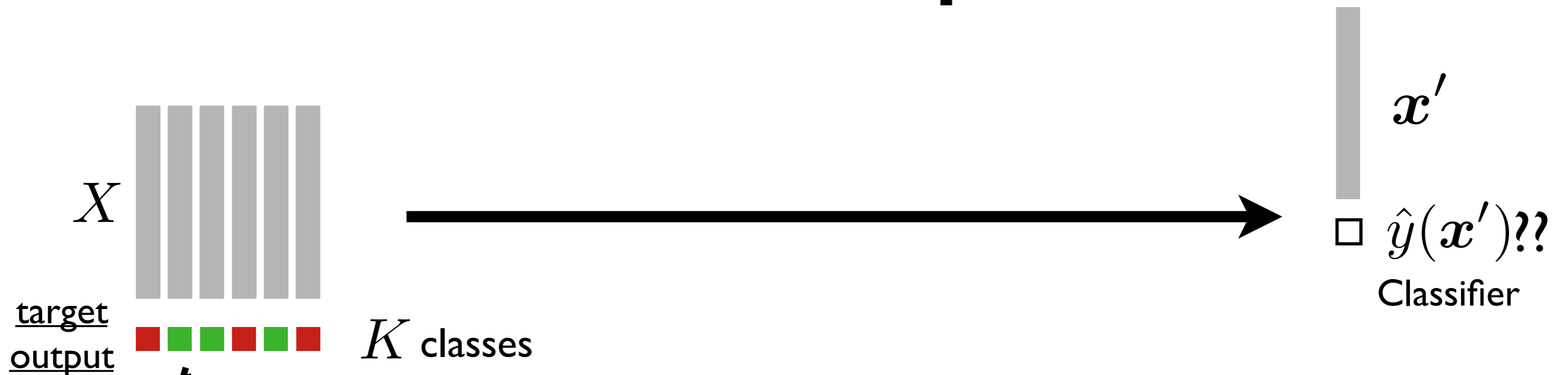
Compressive Classification
(a proof of concept)



Classification problem

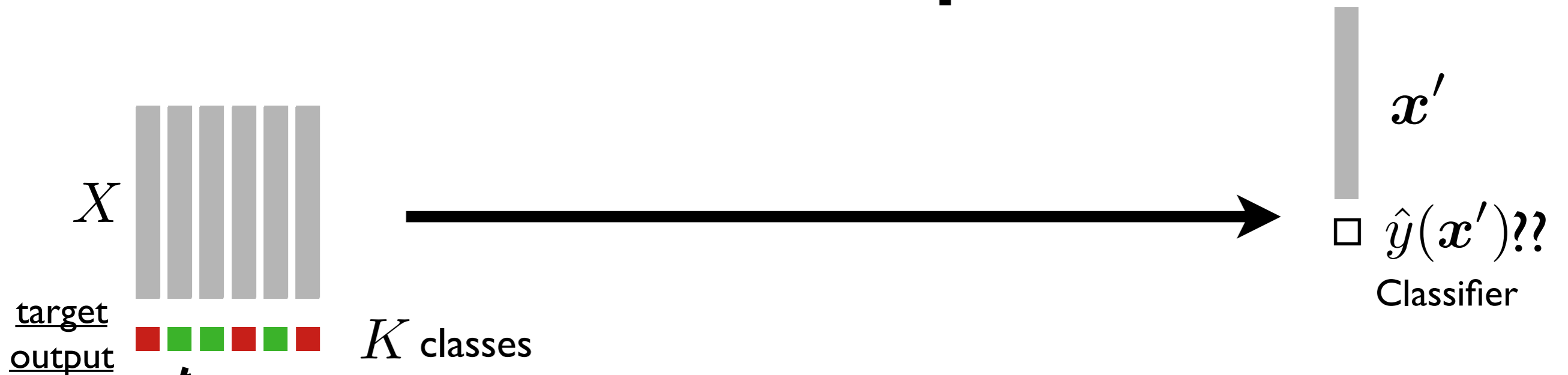


Classification problem



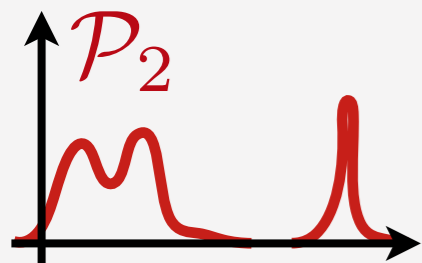
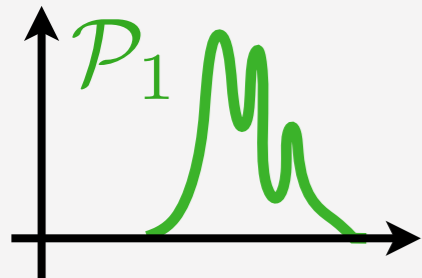
Model: $(\mathbf{x}_i, y_i) \sim \sum_{k=1}^K p_k \mathcal{P}_k(\mathbf{x})$

Classification problem

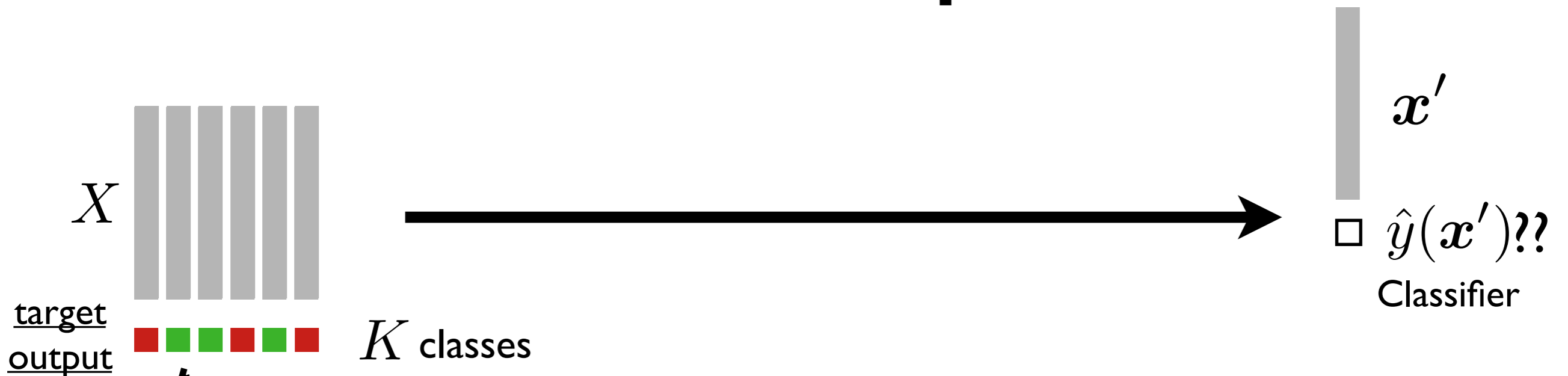


Model: $(\mathbf{x}_i, y_i) \sim \sum_{k=1}^K p_k \mathcal{P}_k(\mathbf{x})$

Conditional
 $\mathcal{P}(\mathbf{x} | y = k)$

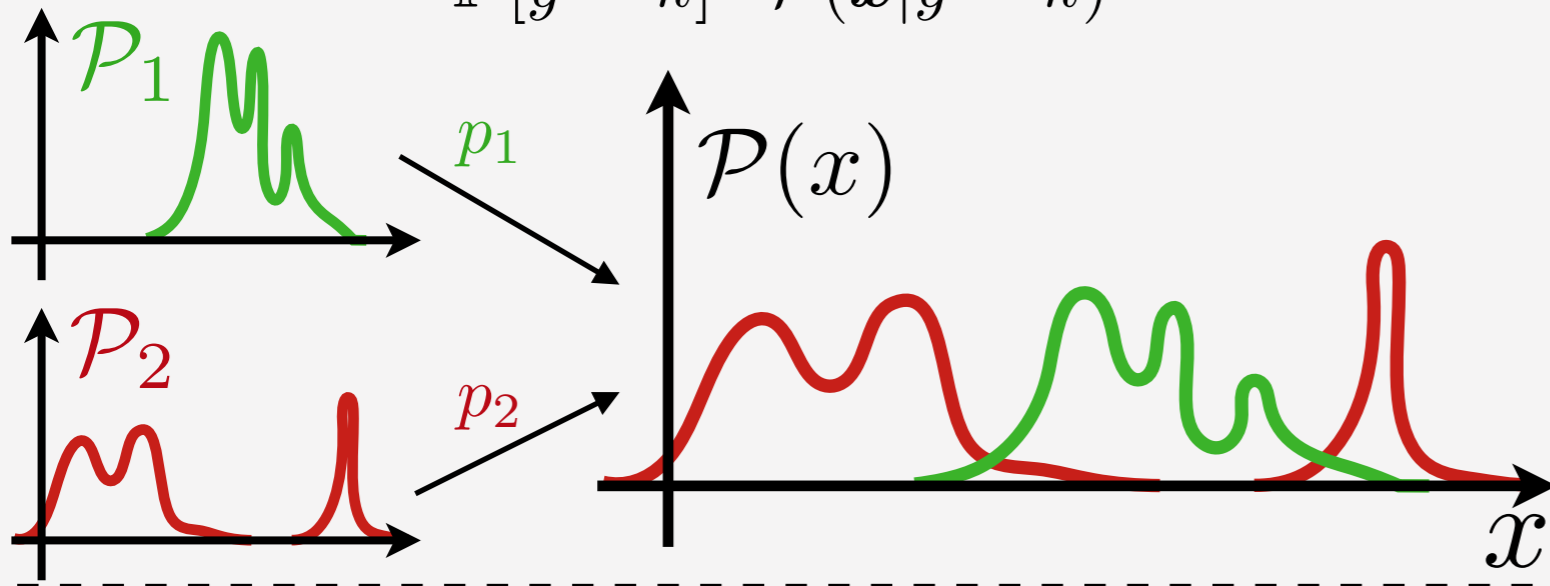


Classification problem

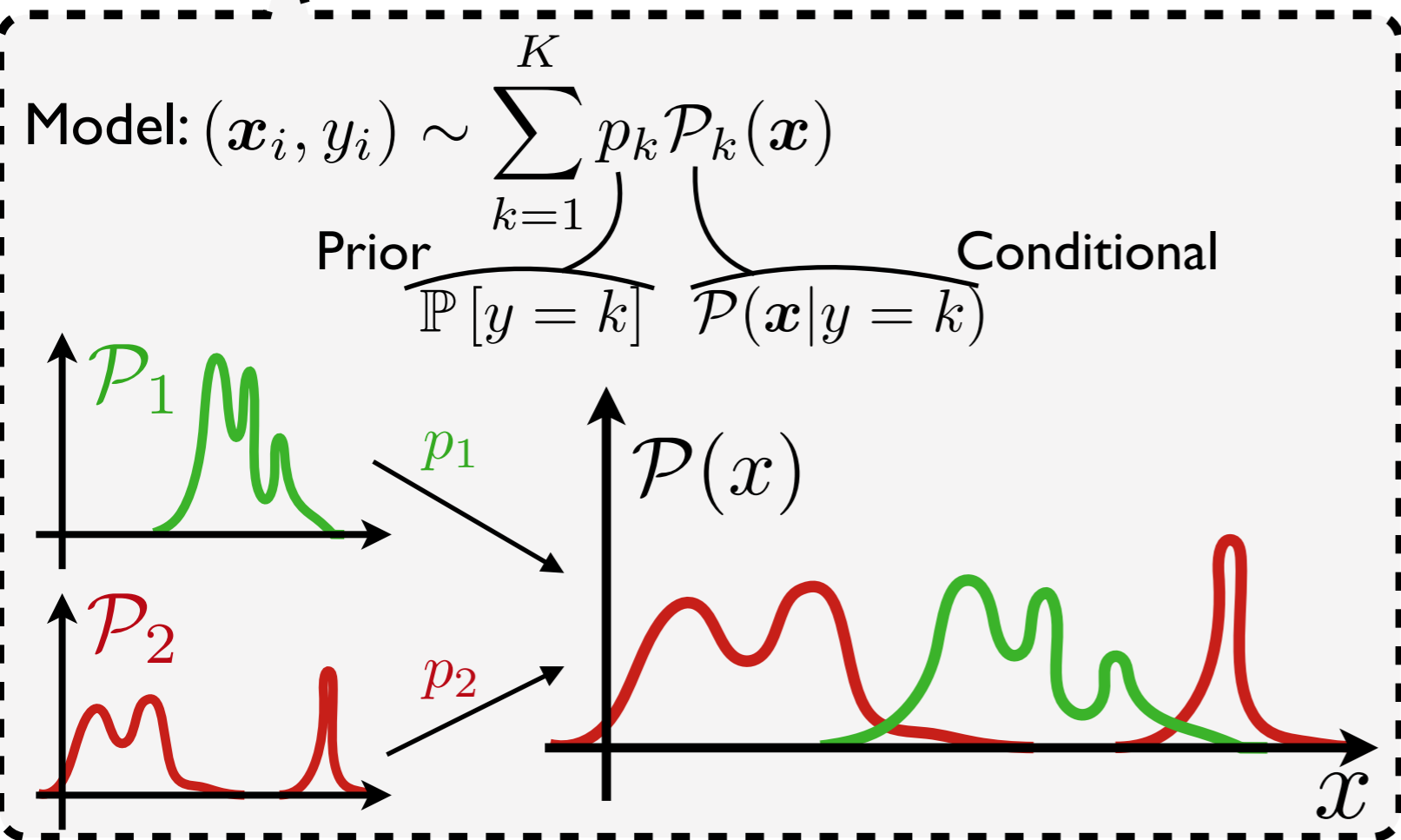
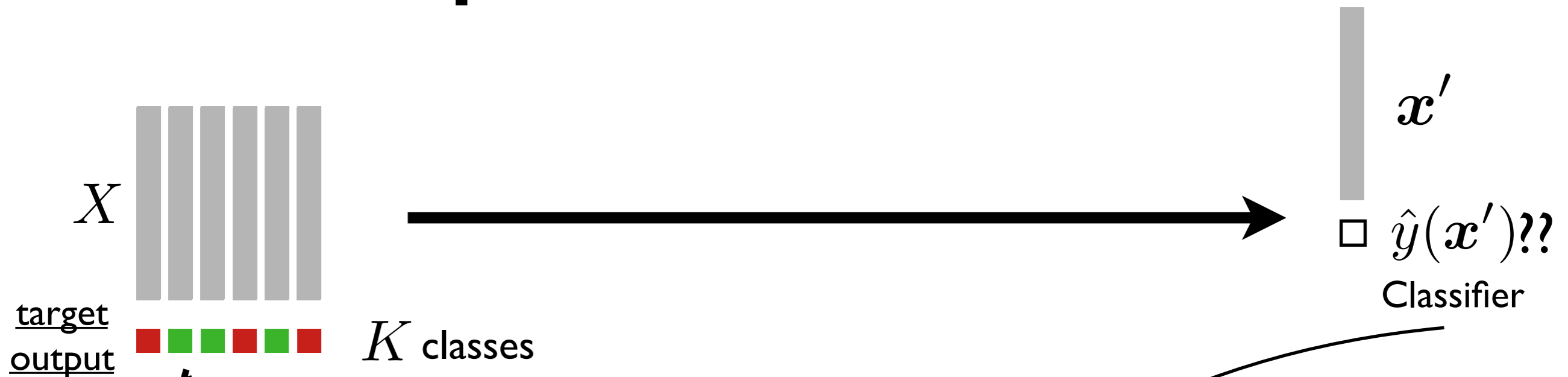


Model: $(\mathbf{x}_i, y_i) \sim \sum_{k=1}^K p_k \mathcal{P}_k(\mathbf{x})$

Prior $\mathbb{P}[y = k]$ Conditional $\mathcal{P}(\mathbf{x}|y = k)$

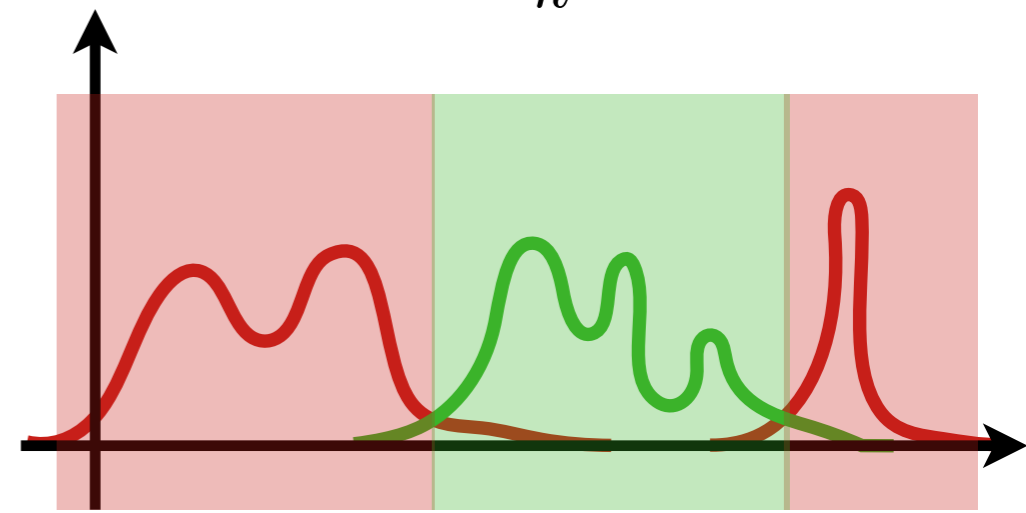


Optimal classifier

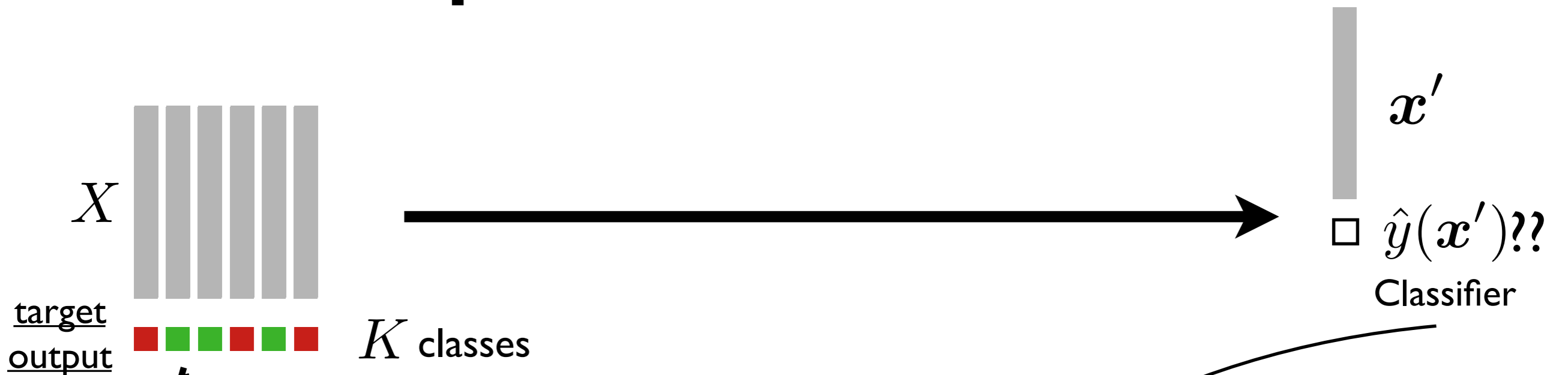


Best we can do (oracle knowledge):
MAP (or Bayes) classifier

$$\hat{y}^{\text{MAP}} := \arg \max_k p_k \mathcal{P}_k(\mathbf{x}')$$

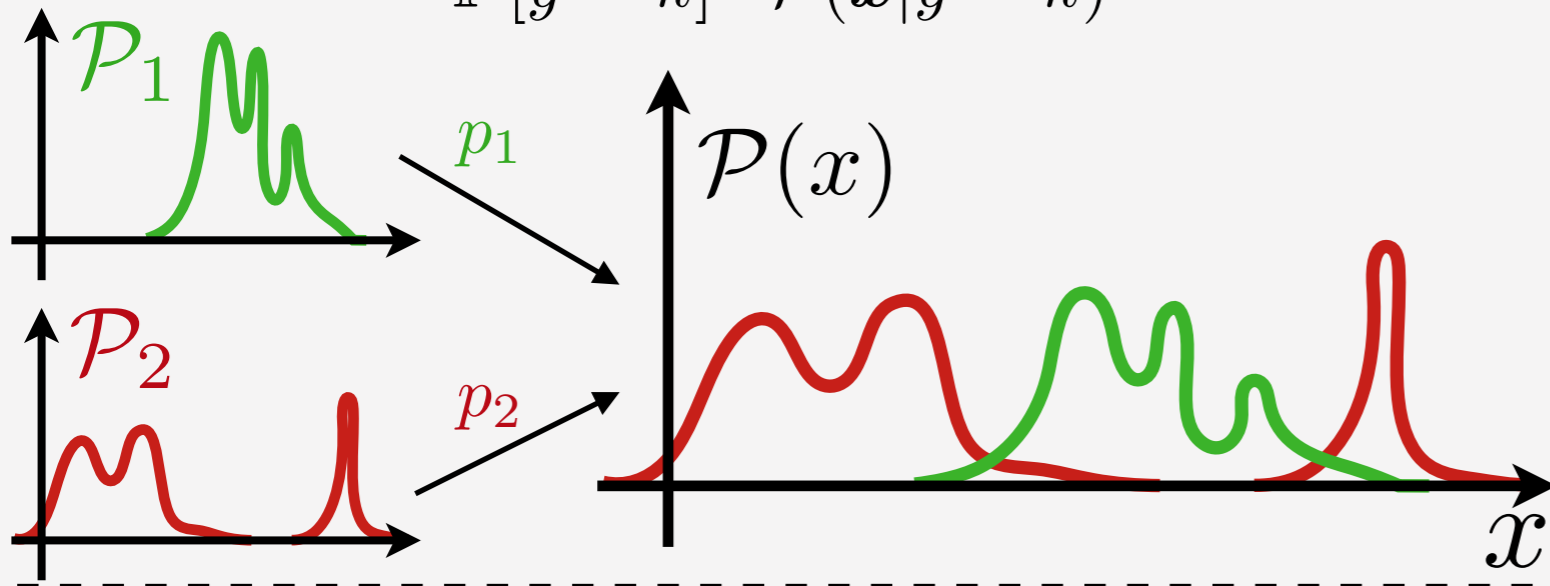


Optimal classifier



Model: $(\mathbf{x}_i, y_i) \sim \sum_{k=1}^K p_k \mathcal{P}_k(\mathbf{x})$

Prior $\mathbb{P}[y = k]$ Conditional $\mathcal{P}(\mathbf{x}|y = k)$

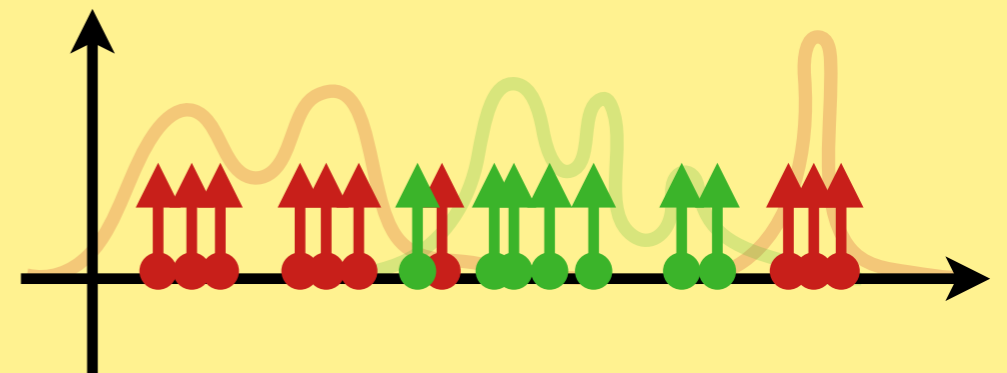


Best we can do (oracle knowledge):
MAP (or Bayes) classifier

$$\hat{y}^{\text{MAP}} := \arg \max_k p_k \mathcal{P}_k(\mathbf{x}')$$

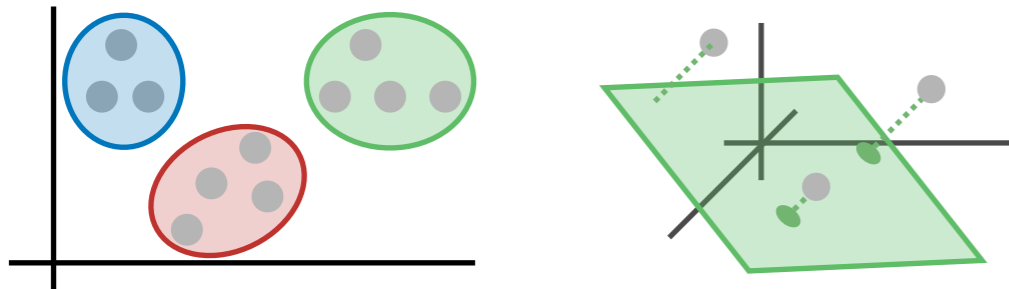
Hard to estimate in practice!

You get this:



In this talk... (finally)

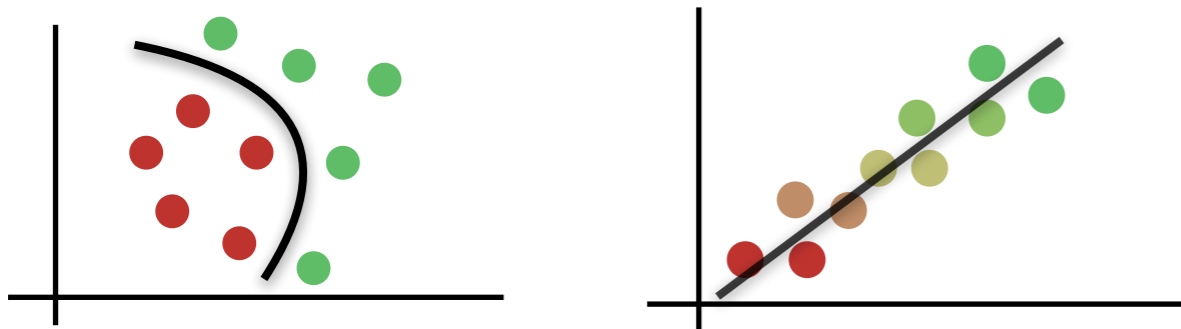
Unsupervised ML



Unsupervised Compressive Learning

- Compressive K-Means [Keriven-CKM]
- Compressive GMM estimation [Keriven-GMM]
- Compressive PCA [Gribonval-CL]

Supervised ML



Supervised Compressive Learning

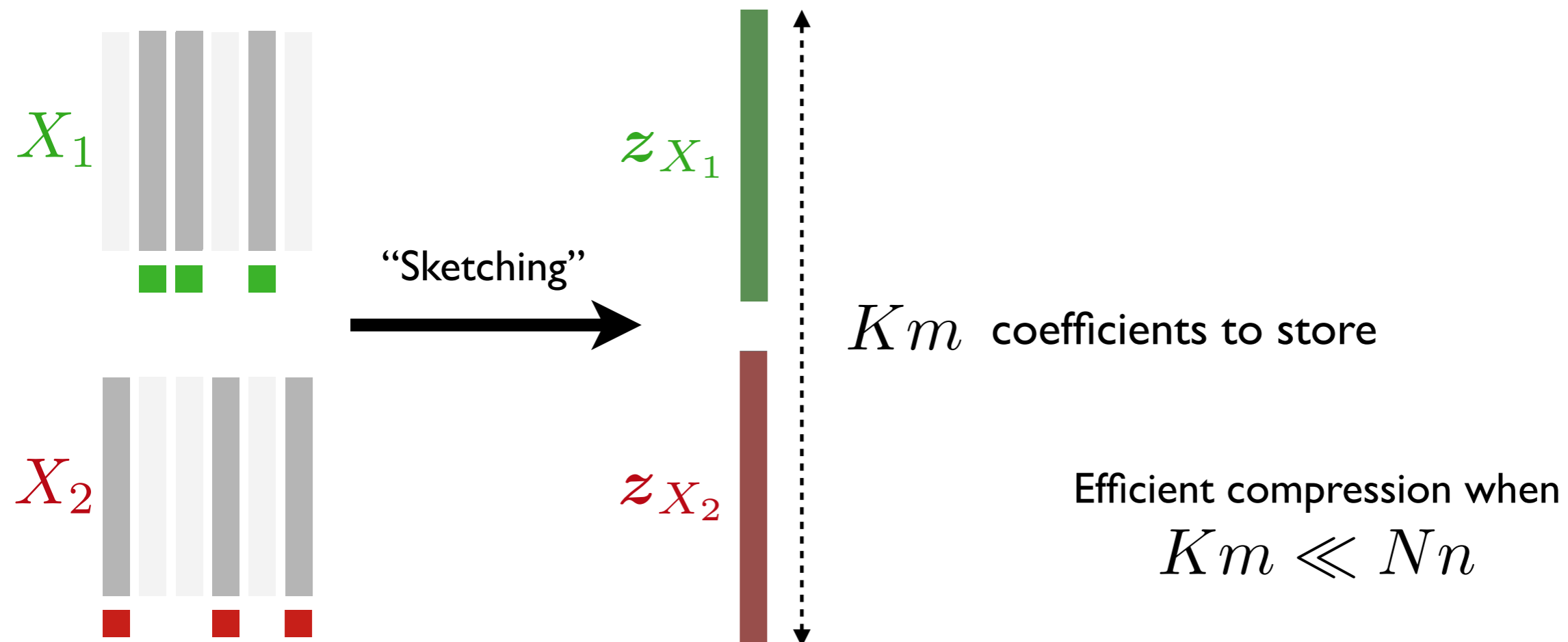
Compressive Classification
(a proof of concept)

Compressive Classifier

Observation (sketching) phase

One sketch per class! $z_{X_k} = \frac{1}{N_k} \sum_{\mathbf{x}_i \in X_k} f(\mathbf{x}_i)$
 $k \in \{1, \dots, K\} \quad \simeq$

i.e., sketch the conditionals! $\mathcal{A}(\mathcal{P}_k) = \mathbb{E}_{\mathbf{x} \sim \mathcal{P}_k} f(\mathbf{x})$

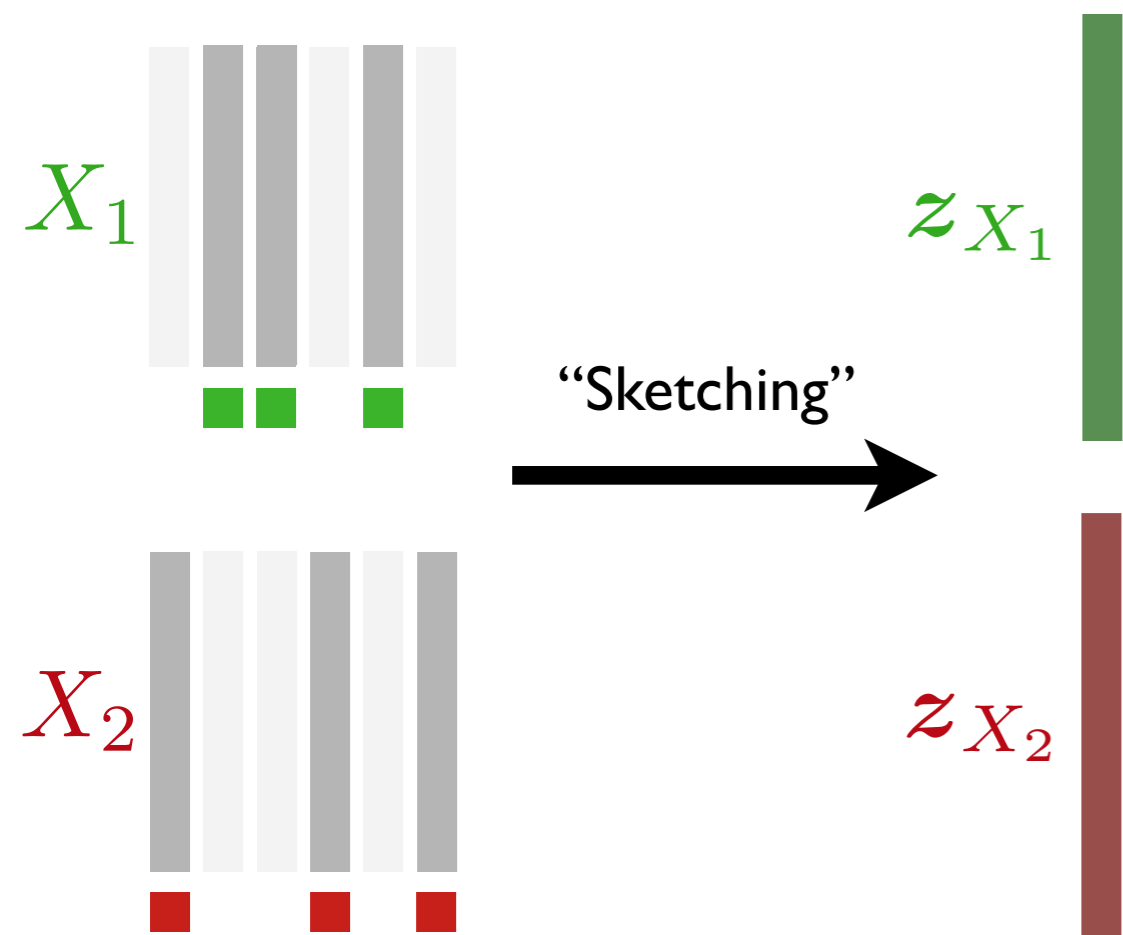


Compressive Classifier

Observation (sketching) phase

One sketch per class! $z_{X_k} = \frac{1}{N_k} \sum_{\mathbf{x}_i \in X_k} f(\mathbf{x}_i)$
 $k \in \{1, \dots, K\} \quad \simeq$

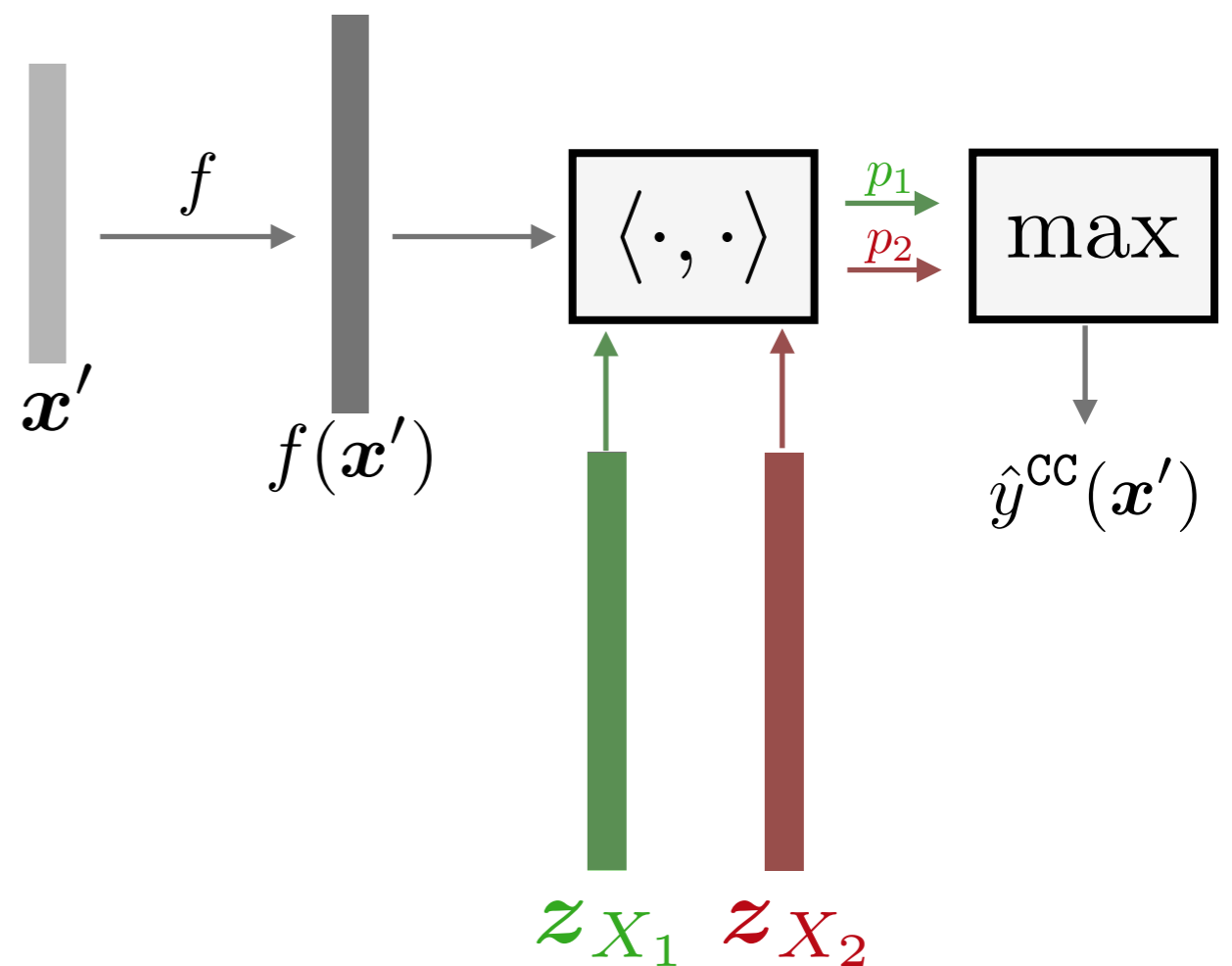
i.e., sketch the $\mathcal{A}(\mathcal{P}_k) = \mathbb{E}_{\mathbf{x} \sim \mathcal{P}_k} f(\mathbf{x})$
 conditionals!



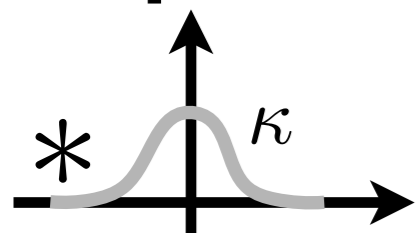
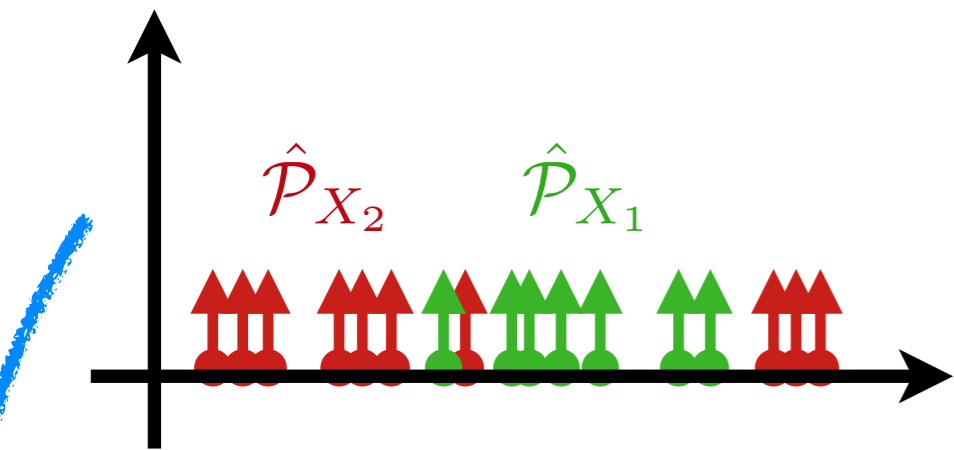
Classification phase

$\hat{y}^{\text{CC}}(\mathbf{x}') := \arg \max_k \hat{p}_k \cdot \langle f(\mathbf{x}'), z_{X_k} \rangle$

$\frac{N_k}{N}$ approximated prior sketch class correlation



Interpretation



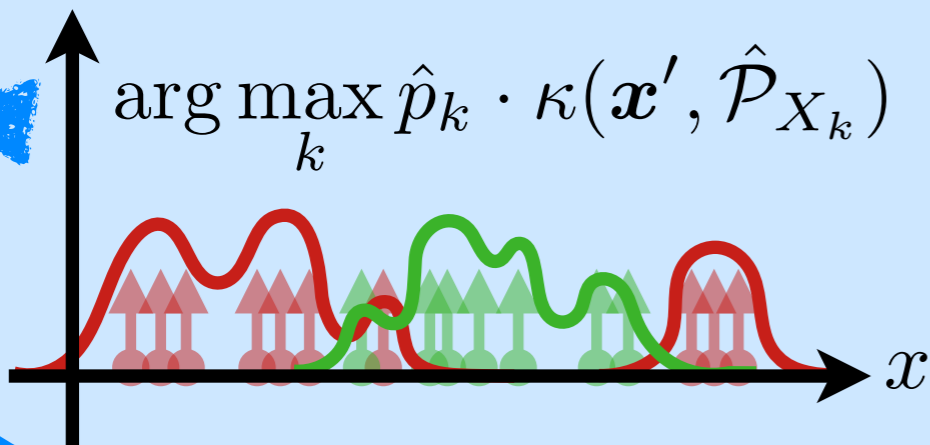
$$\langle f(\mathbf{x}), f(\mathbf{x}') \rangle \simeq \kappa(\mathbf{x}, \mathbf{x}')$$

$$\kappa(\cdot, \mathcal{P}) := \mathbb{E}_{\mathbf{x}' \sim \mathcal{P}} \kappa(\cdot, \mathbf{x}')$$

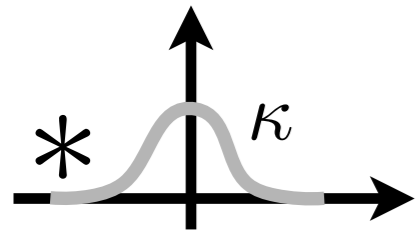
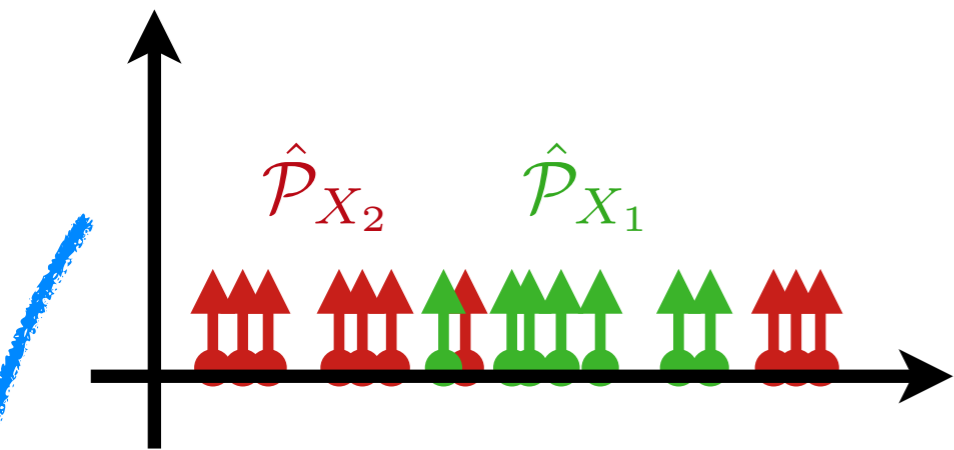
$$m \rightarrow \infty$$

\mathcal{H}_κ

$$\arg \max_k \hat{p}_k \cdot \kappa(\mathbf{x}', \hat{\mathcal{P}}_{X_k})$$



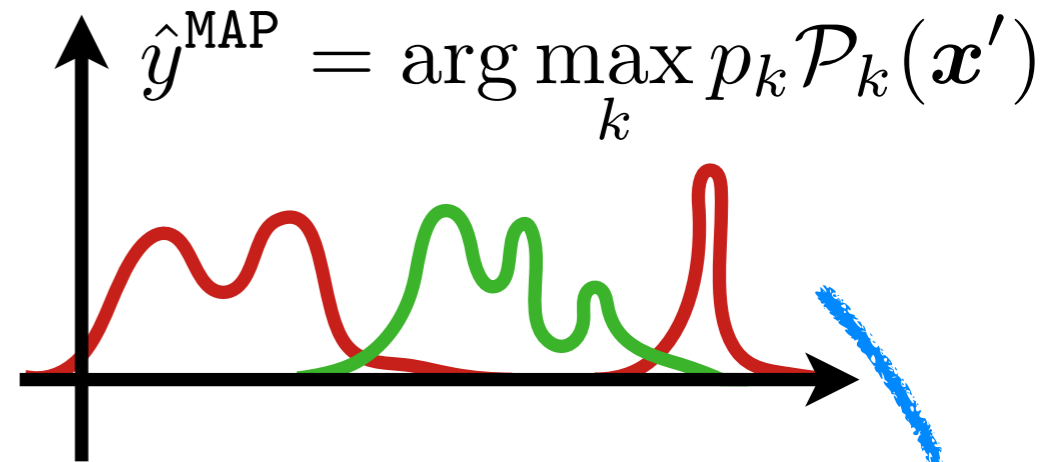
Interpretation



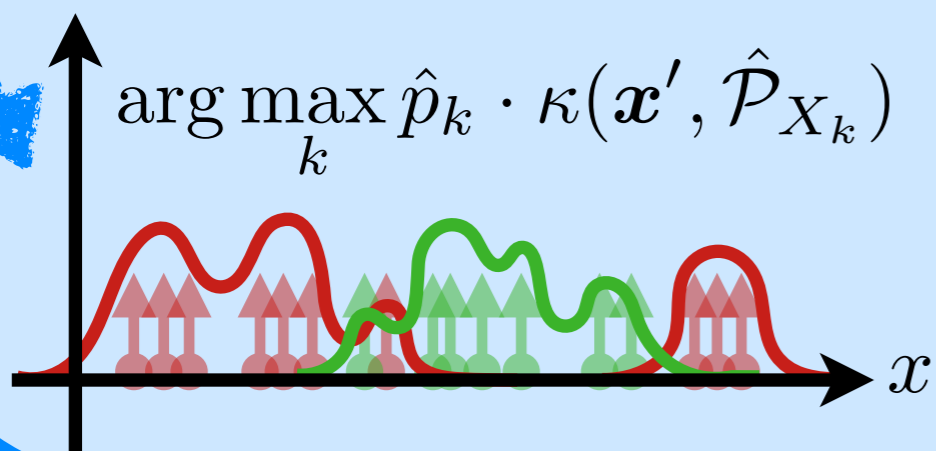
$$\langle f(\mathbf{x}), f(\mathbf{x}') \rangle \simeq \kappa(\mathbf{x}, \mathbf{x}')$$

$$\kappa(\cdot, \mathcal{P}) := \mathbb{E}_{\mathbf{x}' \sim \mathcal{P}} \kappa(\cdot, \mathbf{x}')$$

$$m \rightarrow \infty$$

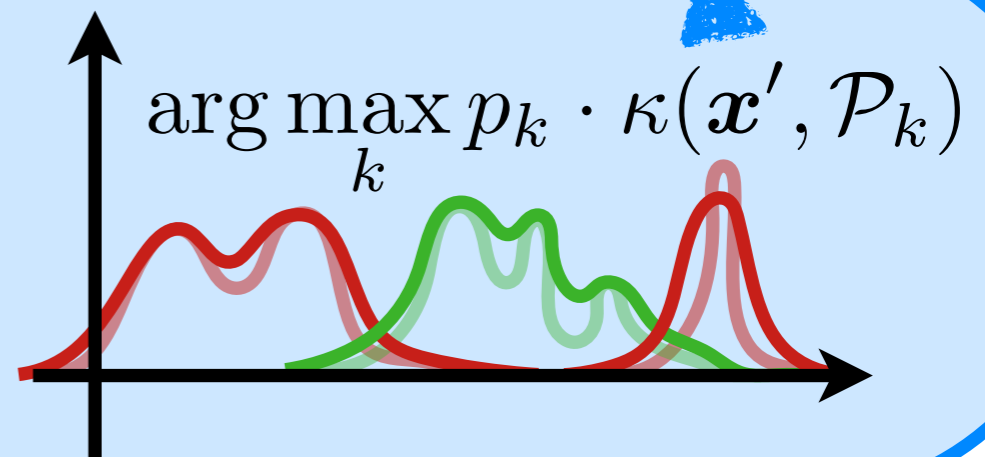


\mathcal{H}_κ

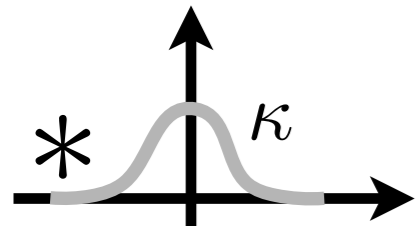
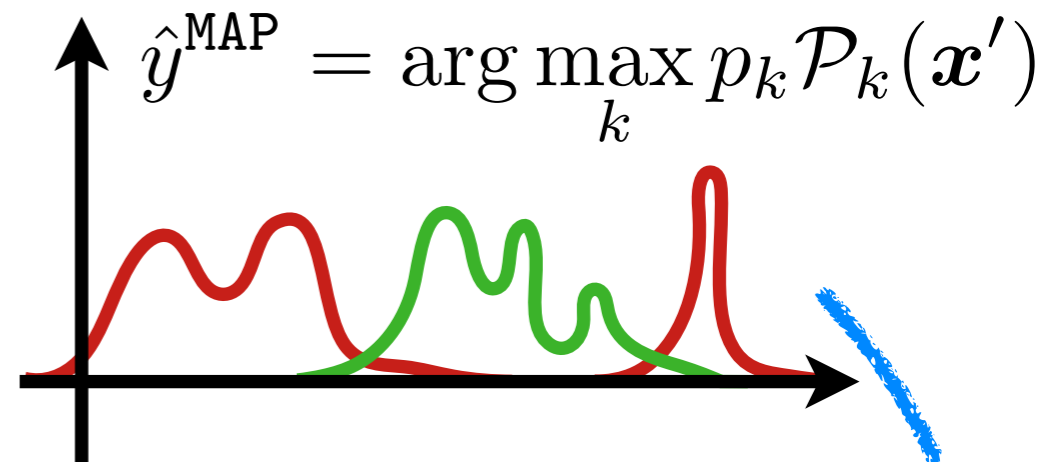
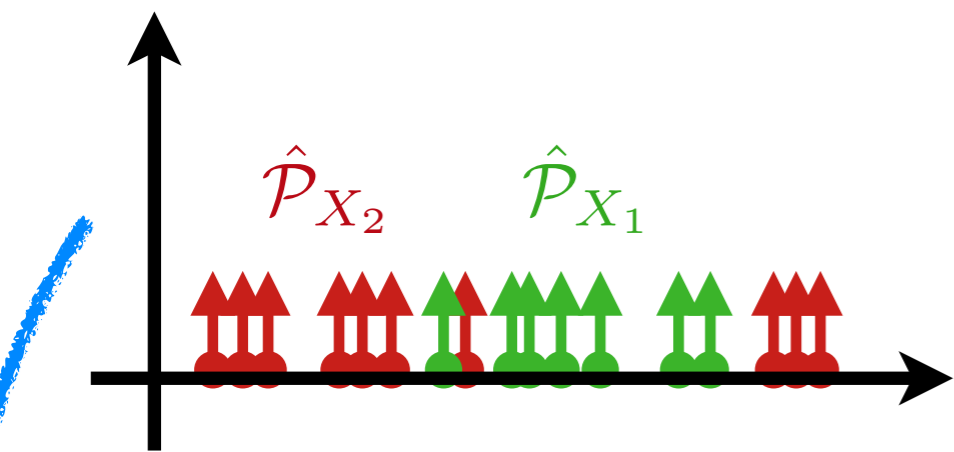


$$N \rightarrow \infty$$

$$\simeq$$



Interpretation



CC approximates the MAP classifier
projected into the RKHS \mathcal{H}_κ !

$$\langle f(\mathbf{x}), f(\mathbf{x}') \rangle \simeq \kappa(\mathbf{x}, \mathbf{x}')$$

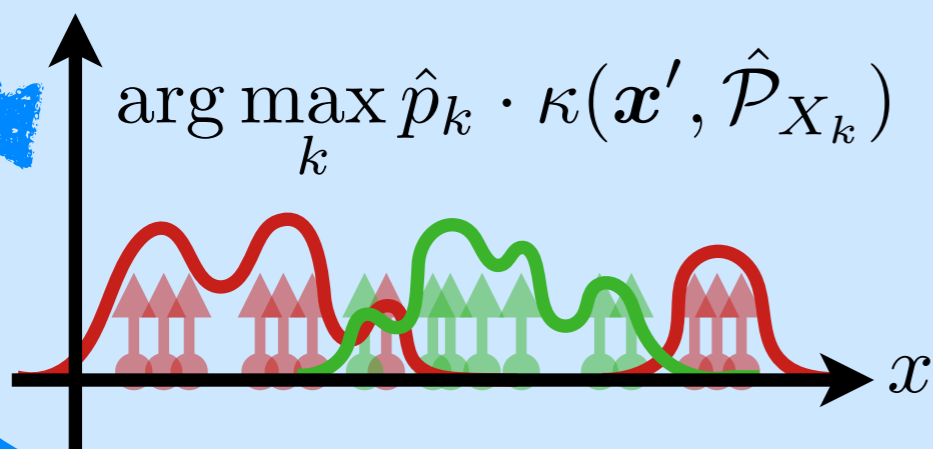
$$\kappa(\cdot, \mathcal{P}) := \mathbb{E}_{\mathbf{x}' \sim \mathcal{P}} \kappa(\cdot, \mathbf{x}')$$

$$m \rightarrow \infty$$

$$\hat{y}^{\text{CC}} = \arg \max_k \hat{p}_k \cdot \langle f(\mathbf{x}'), \mathbf{z}_{X_k} \rangle$$

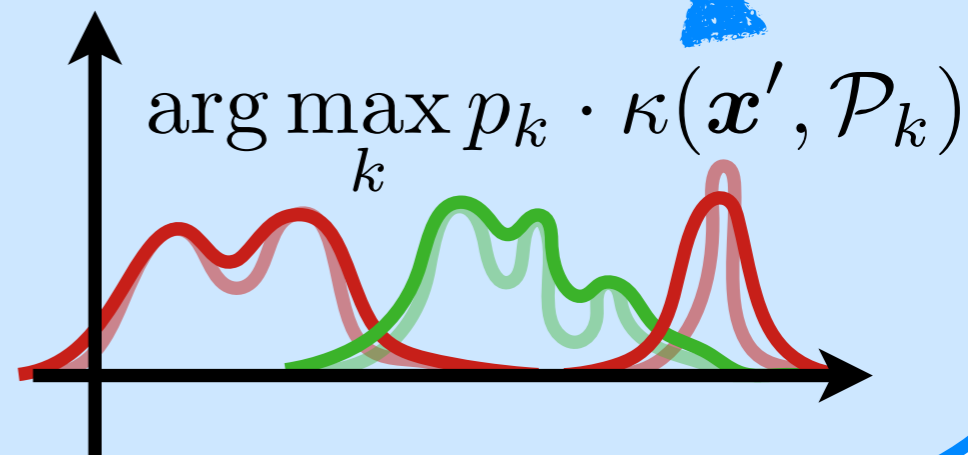
$$\simeq \arg \max_k p_k \cdot \kappa(\mathbf{x}', \mathcal{P}_k)$$

\mathcal{H}_κ



$$N \rightarrow \infty$$

$$\simeq$$



Pros and cons

$$z_{X_k} = \frac{1}{N_k} \sum_{\mathbf{x}_i \in X_k} f(\mathbf{x}_i) \quad \text{then} \quad \hat{y}^{\text{cc}} = \arg \max_k \hat{p}_k \cdot \langle f(\mathbf{x}'), z_{X_k} \rangle$$
$$\simeq \arg \max_k p_k \cdot \kappa(\mathbf{x}', \mathcal{P}_k)$$

Pro

- Cheap to “learn”: you only observe the dataset (“compressive k-NN”)
- Cheap to evaluate: f (once), $+$, $*$ and \max
- Easy to parallelize/update: suited to massive datasets, distributed (sensitive?) datasets, data streams, and data augmentation (almost free!)
- MAP in RKHS interpretation is seducing, but it is an asymptotical result...

Pros and cons

$$z_{X_k} = \frac{1}{N_k} \sum_{\mathbf{x}_i \in X_k} f(\mathbf{x}_i) \quad \text{then} \quad \hat{y}^{\text{CC}} = \arg \max_k \hat{p}_k \cdot \langle f(\mathbf{x}'), z_{X_k} \rangle$$
$$\simeq \arg \max_k p_k \cdot \kappa(\mathbf{x}', \mathcal{P}_k)$$

Pro

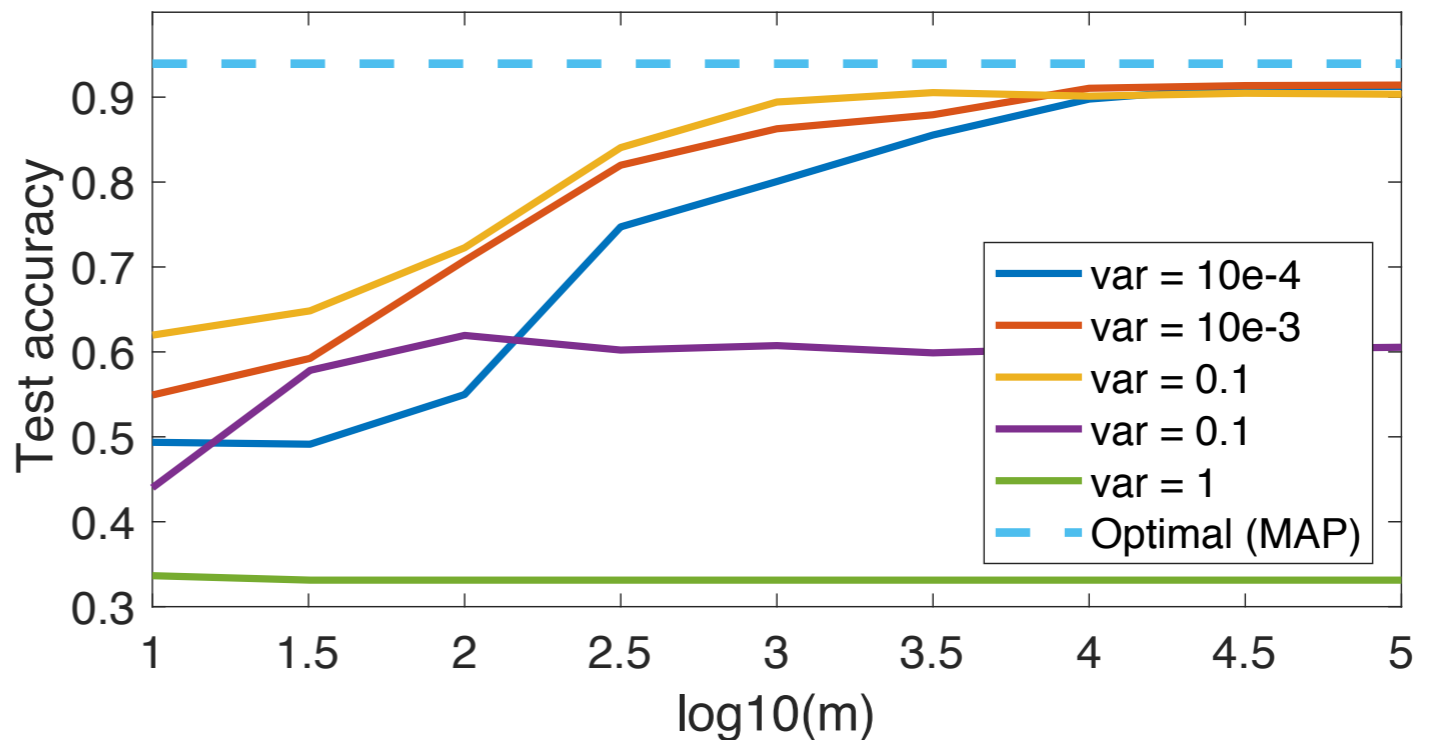
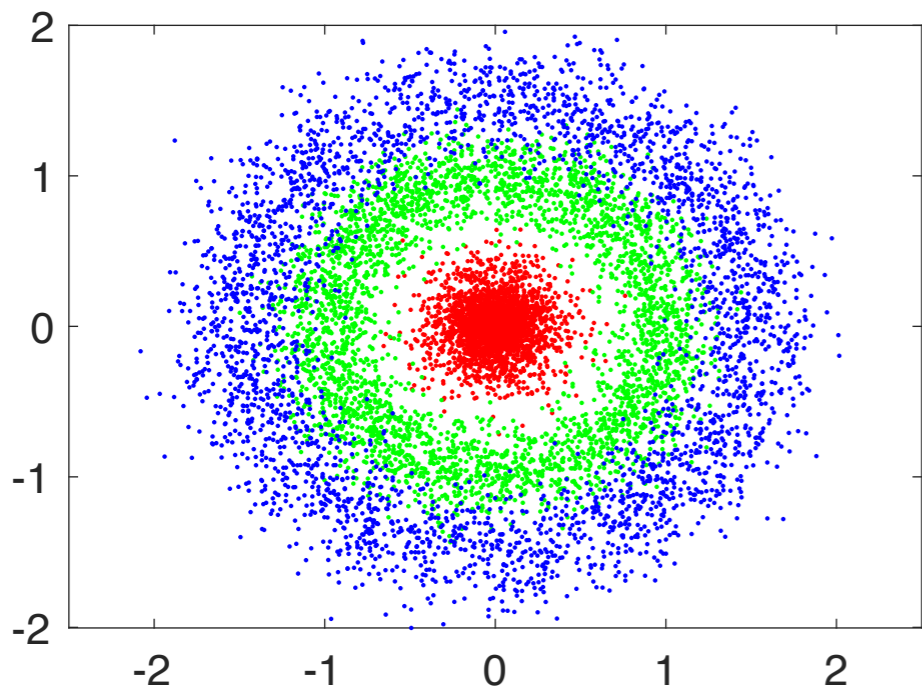
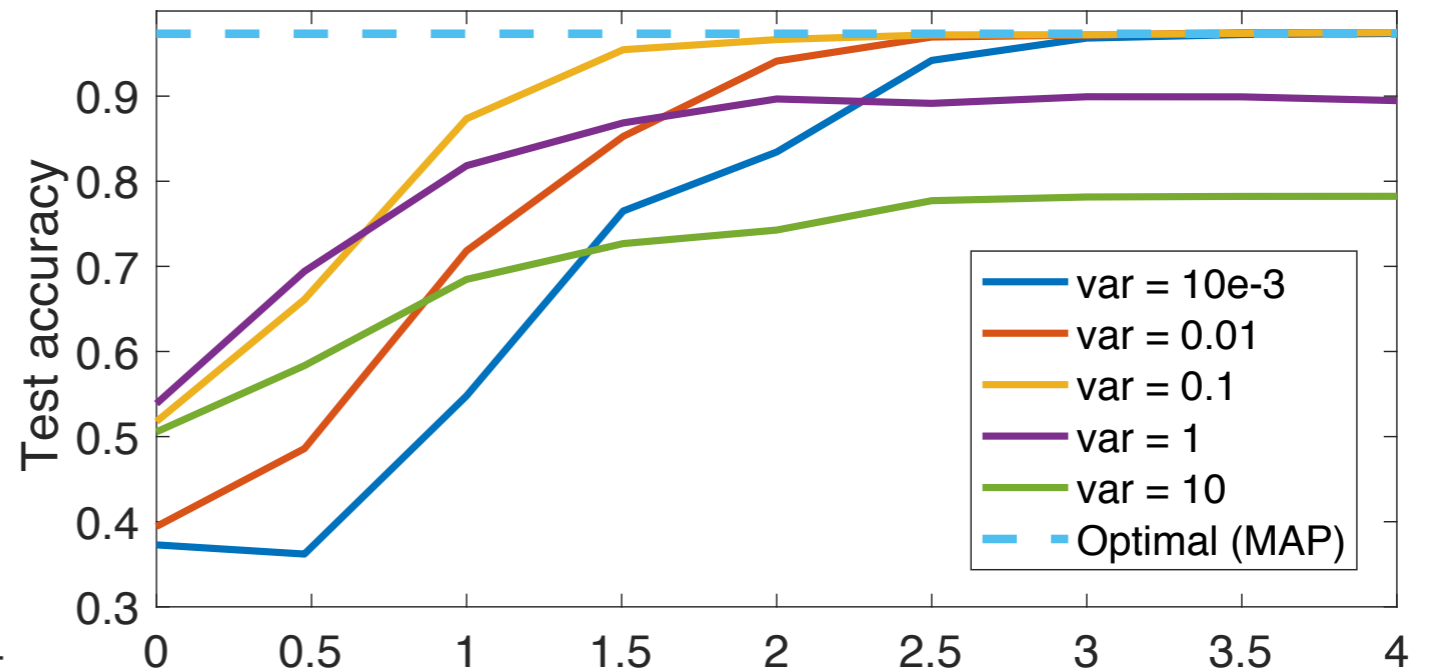
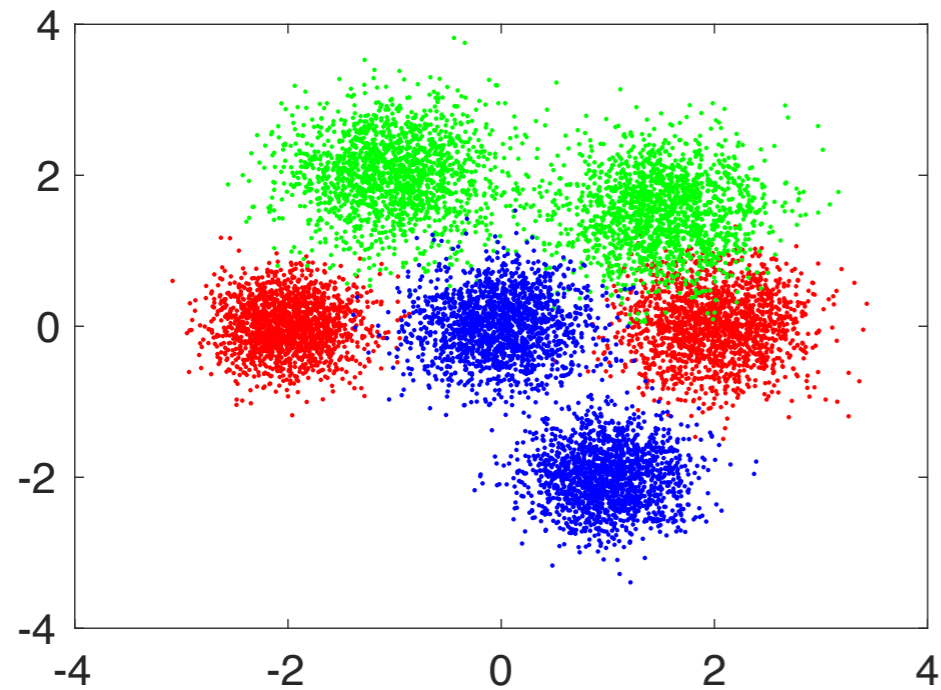
- Cheap to “learn”: you only observe the dataset (“compressive k-NN”)
- Cheap to evaluate: f (once), $+$, $*$ and \max
- Easy to parallelize/update: suited to massive datasets, distributed (sensitive?) datasets, data streams, and data augmentation (almost free!)
- MAP in RKHS interpretation is seducing, but it is an asymptotical result...

Con

- No guarantees (yet): generalization error bounds, sample complexity bounds...
- Hyperparameters: choice of f (hence κ) is crucial!
- High-Dimensional data is expected to be difficult

Proof of concept (synthetical)

Random Fourier Features sketch with Gaussian kernel




Proof of concept (real)


Random Fourier Features sketch with Gaussian kernel

Standard learned classifier

Compressive classifier



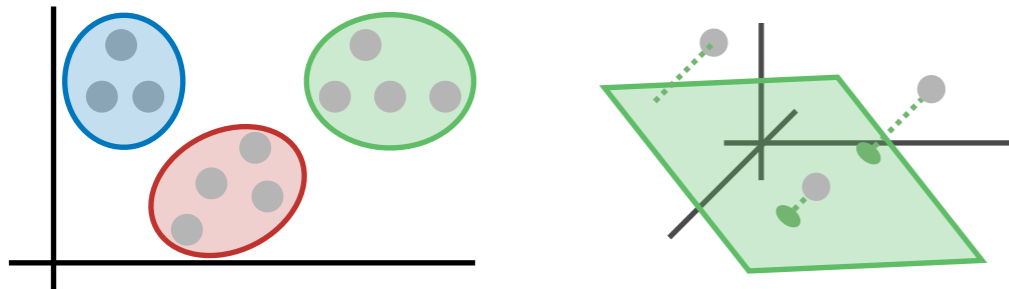
	N	n	K	SVM	$m = 50$	$m = 1000$
Iris	150	4	3	2.00	6.51 ± 1.81	5.51 ± 1.23
				4.00	8.22 ± 3.25	6.18 ± 2.40
Wine	178	13	3	0.84	4.56 ± 2.34	2.43 ± 0.72
				1.69	13.75 ± 4.09	8.19 ± 1.29
Breast cancer	569	30	2	3.67	7.00 ± 1.40	3.93 ± 0.39
				2.13	9.22 ± 2.33	6.23 ± 0.69
Adult (3 attr.)	30718	3	2	21.03	23.88 ± 4.37	23.11 ± 1.05
				21.06	36.09 ± 6.67	35.04 ± 1.63



Error rates (%) with SD
on train and test sets

In this talk...

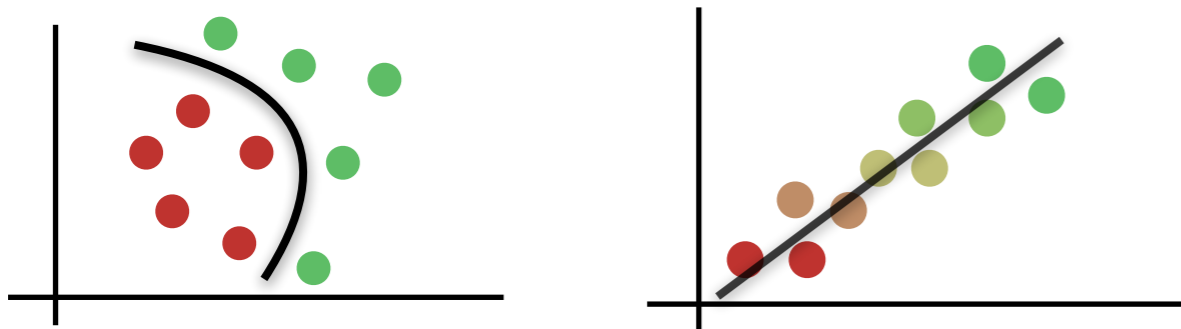
Unsupervised ML



Unsupervised Compressive Learning

- Compressive K-Means [Keriven-CKM]
- Compressive GMM estimation [Keriven-GMM]
- Compressive PCA [Gribonval-CL]

Supervised ML



Supervised Compressive Learning


Compressive Classification
(a proof of concept)

What about HD data such as images?


Sketching images

Up to now, the sketch function f used are Random Fourier Features

$$f(\mathbf{x})_{\text{RFF}} = \left[\exp(i\omega_j^T \mathbf{x}) \right]_{j=1}^m$$




$$\kappa(\mathbf{x}, \mathbf{x}') = K(\mathbf{x} - \mathbf{x}') \quad \text{Depends only on signal difference}$$



For images, a kernel based on the pixel-wise difference does not make so much sense...

Idea: capture the image structure with a (random) Convolutional Neural Network architecture

$$f(\mathbf{x})_{\text{CNN}} = \text{CNN}_{\theta}(\mathbf{x})$$



Random CNN have been shown to be surprisingly meaningful features for image processing tasks
E.g., [Cho], [Giryas], [Rosenfeld], ...

Sketching images

Up to now, the sketch function f used are Random Fourier Features


$$f(\mathbf{x})_{\text{RFF}} = \left[\exp(i\omega_j^T \mathbf{x}) \right]_{j=1}^m$$

$$\Updownarrow$$

$$\kappa(\mathbf{x}, \mathbf{x}') = K(\mathbf{x} - \mathbf{x}')$$

For images, a kernel based on the pixel-wise difference does not make much sense...

Idea: capture the image structure with a (random) Convolutional Neural Network architecture

$$f(\mathbf{x})_{\text{CNN}} = \text{CNN}_{\theta}(\mathbf{x})$$


	N	n	CNN	m = 250	m = 5000
MNIST	60000	28 × 28 × 1	1.60 ± 0.12	17.73 ± 1.43	16.60 ± 1.54
	10000		1.63 ± 0.11	16.83 ± 1.39	15.80 ± 1.61
CIFAR10	50000	32 × 32 × 3	39.08 ± 1.48	71.76 ± 1.85	72.83 ± 2.00
	10000		40.28 ± 1.36	71.12 ± 1.72	72.02 ± 1.85

“Not that random” (Work In Progress)...

Some perspectives

$$z_{X_k} = \frac{1}{N_k} \sum_{\mathbf{x}_i \in X_k} f(\mathbf{x}_i) \quad \text{then} \quad \hat{y}^{\text{cc}} = \arg \max_k \hat{p}_k \cdot \langle f(\mathbf{x}'), z_{X_k} \rangle$$
$$\simeq \arg \max_k p_k \cdot \kappa(\mathbf{x}', \mathcal{P}_k)$$

- Explore the features/kernel choice?
- Learn the kernel from a small data sample (distilled sensing)?
- (Non-asymptotical) formal guarantees?
- Is this the best we can do with the sketch?
- Regression? Other supervised tasks?
- ...

(Some) references

- [Keriven-CKM] N. Keriven, N. Tremblay, Y. Traonmilin, R. Gribonval, “Compressive K-Means”, ICASSP 2017
- [Keriven-GMM] N. Keriven, A. Bourrier, R. Gribonval, P. Perez, “Sketching for Large-Scale Learning of Mixture Models”, Information and Inference, 2017
- [Gribonval-CL] R. Gribonval, G. Blanchard, N. Keriven, Y. Traonmilin, “Compressive Statistical Learning with Random Feature Moments”, Arxiv, 2017
- [Giryes] R. Giryes, G. Sapiro, A.M. Bronstein, “Deep Neural Networks with Random Gaussian Weights: A Universal Classification Strategy?”, IEEE Transactions on Signal Processing, 2016
- [Cho] Y. Cho, L.K. Saul, “Kernel methods for deep learning”, NIPS 2009
- [Rosenfeld] A. Rosenfeld, J.K. Tsotsos, “Intriguing Properties of Randomly Weighted Networks: Generalizing While Learning Next to Nothing”, Arxiv, 2018
- [Smola] A. Smola, A. Gretton, L. Song, B. Scholkopf, “A Hilbert space embedding for distributions”, International Conference on Algorithmic Theory, 2007