



FSAN/ELEG815: Statistical Learning

Gonzalo R. Arce

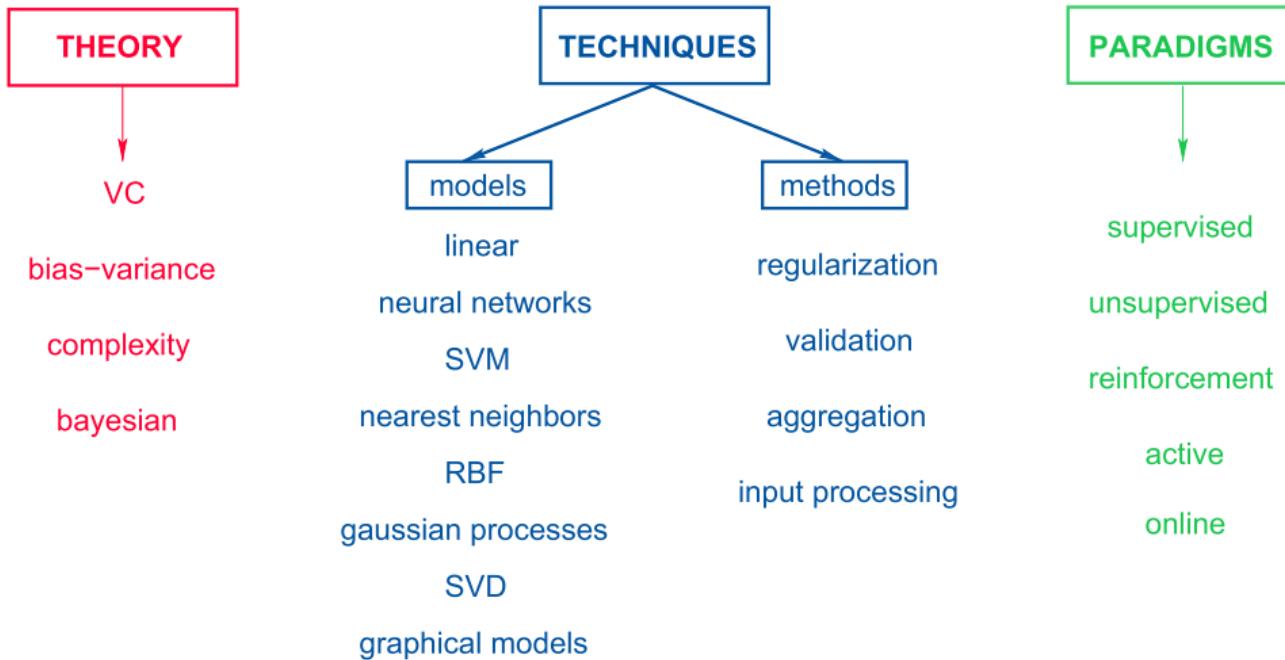
Department of Electrical and Computer Engineering
University of Delaware

Epilogue

It is a Jungle Out There

semi-supervised learning **overfitting** stochastic gradient descent **SVM** *Q* learning
Gaussian processes **deterministic noise** data snooping learning curves
distribution-free linear regression VC dimension mixture of expe
collaborative filtering nonlinear transformation **sampling bias** neural networks
decision trees RBF training versus testing noisy targets no free
active learning linear models bias-variance tradeoff Bayesian prior
ordinal regression cross validation logistic regression data contamination
ensemble learning error measures types of learning perceptrons hidden Markov mo
exploration versus exploitation **is learning feasible?** kernel methods graphical models
clustering regularization weight decay soft-order constraint
Boltzmann mach

The Map

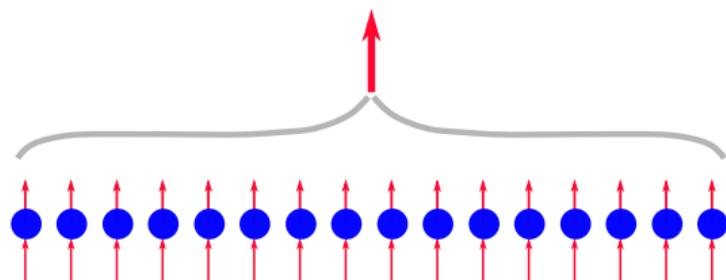


Outline

1. Aggregation.
2. Kernel Methods.
3. Randomized Singular Value Decomposition.

What is Aggregation?

Combining different solutions h_1, h_2, \dots, h_T that were trained on \mathcal{D}

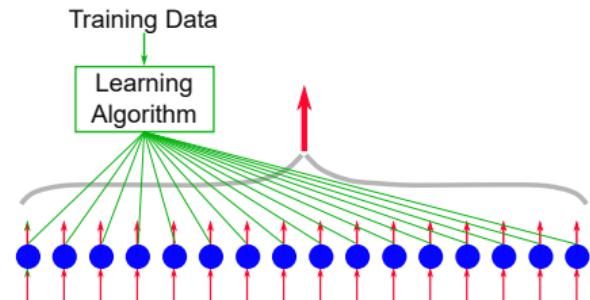


Regression: take an average.
Classification: take a vote.

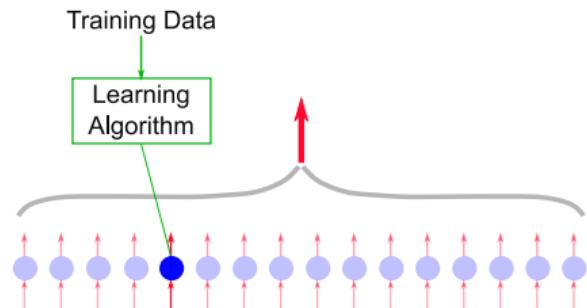
a.k.a **ensemble learning** and **boosting**.

Different from 2-Layer Learning

In a 2-layer model, all units learn **jointly**:



In aggregation, they learn **independently** then get combined:



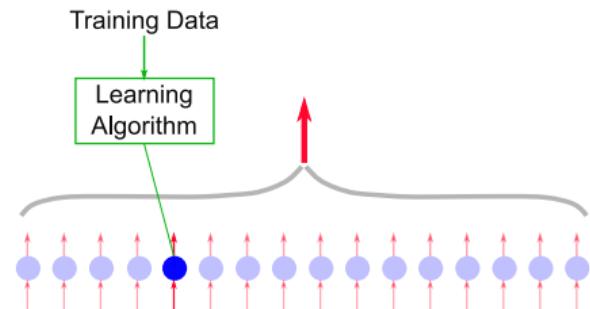
Two Types of Aggregation

After the fact: combines existing solutions

Example: Netflix teams merging "blending".

Before the fact: creates solutions to be combined

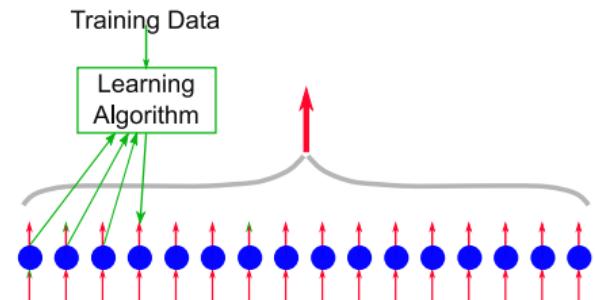
Example: Bagging - resampling \mathcal{D} .



Decorrelation - Boosting

Create h_1, \dots, h_t, \dots sequentially:

Makes h_t decorrelated with previous h 's:



Emphasize points in \mathcal{D} that were misclassified.
Choose weight of h_t based on $E_{\text{in}}(h_t)$.

Blending - After the Fact

For regression, $h_1, h_2, \dots, h_T \rightarrow g(x) = \sum_{t=1}^T \alpha_t h_t(\mathbf{x})$

Principled choice of α_t 's: minimize the error on an "aggregation data set"

pseudo-inverse

Some α_t 's can come out negative.

Most valuable h_t in the blend?

Outline

1. Aggregation.
2. Kernel Methods.
3. Randomized Singular Value Decomposition.

Support Vector Machines - "Support Vectors" in \mathcal{X} Space

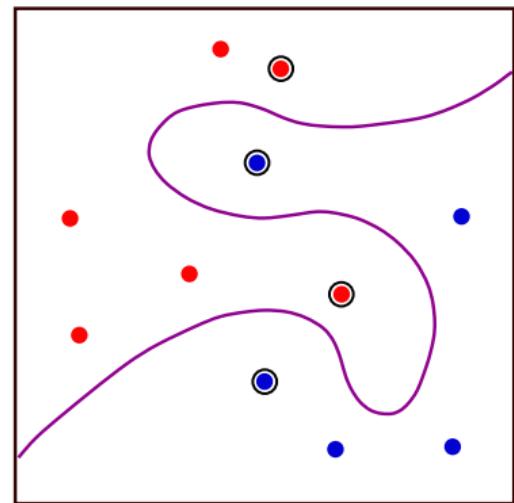
Support vectors live in the \mathcal{Z} space.

In the \mathcal{X} space, "pre-images" of support vectors.

The margin is maintained in the \mathcal{Z} space.

Generalization result

$$\mathbb{E}[E_{out}] \leq \frac{\mathbb{E}[\# \text{ of SV's}]}{N-1}$$



What Do We Need from the \mathcal{Z} Space?

$$\mathcal{L}(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m \mathbf{z}_n^\top \mathbf{z}_m$$

Constraints: $\alpha_n \geq 0$ for $n = 1, \dots, N$ and $\sum_{n=1}^N \alpha_n y_n = 0$.

$$g(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{z} + b) \rightarrow \text{need } \mathbf{z}_n^\top \mathbf{z}.$$

where $\mathbf{w} = \sum_{\mathbf{z}_n \text{ is SV}} \alpha_n y_n \mathbf{z}_n$,

and $b : y_m (\mathbf{w}^\top \mathbf{z} + b) = 1 \rightarrow \text{need } \mathbf{z}_n^\top \mathbf{z}_m$

Generalized Inner Product

Given two points \mathbf{x} and $\mathbf{x}' \in \mathcal{X}$, we need $\mathbf{z}^\top \mathbf{z}'$

Let $\mathbf{z}^\top \mathbf{z}' = K(\mathbf{x}, \mathbf{x}') \rightarrow$ The Kernel "Inner product" of \mathbf{x} and \mathbf{x}' .

Example: $\mathbf{x} = (x_1, x_2) \longrightarrow$ 2nd-order Φ .

$$\mathbf{z} = \Phi(\mathbf{x}) = (1, x_1, x_2, x_1^2, x_2^2, x_1 x_2).$$

$$K(\mathbf{x}, \mathbf{x}') = \mathbf{z}^\top \mathbf{z} = 1 + x_1 x'_1 + x_2 x'_2 + x_1^2 x'^2_1 + x_2^2 x'^2_2 + x_1 x'_1 x_2 x'_2$$

The Trick

Can we compute $K(\mathbf{x}, \mathbf{x}')$ **without** transforming \mathbf{x} and \mathbf{x}' ?

Example: Consider

$$\begin{aligned} K(\mathbf{x}, \mathbf{x}') &= (1 + \mathbf{x}^\top \mathbf{x}')^2 \\ &= (1 + x_1 x'_1 + x_2 x'_2)^2 \\ &= 1 + x_1^2 x'^2_1 + x_2^2 x'^2_2 + 2x_1 x'_1 + 2x_2 x'_2 + 2x_1 x'_1 x_2 x'_2 \end{aligned}$$

This is an inner product!

$$(1, x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2)$$

$$(1, x'^2_1, x'^2_2, \sqrt{2}x'_1, \sqrt{2}x'_2, \sqrt{2}x'_1 x'_2)$$

The Polynomial Kernel

$\mathcal{X} = \mathbb{R}^d$ and $\Phi : \mathcal{X} \rightarrow \mathcal{Z}$ is polynomial of order Q .

The "equivalent" kernel:

$$\begin{aligned} K(\mathbf{x}, \mathbf{x}') &= (1 + \mathbf{x}^\top \mathbf{x}')^Q \\ &= (1 + x_1 x'_1 + x_2 x'_2 + \dots + x_d x'_d)^Q \end{aligned}$$

Compare for $d = 10$ and $Q = 100$

Can adjust scale: $K(\mathbf{x}, \mathbf{x}') = (a\mathbf{x}^\top, \mathbf{x}' + b)^Q$

We Only Need \mathcal{Z} to Exist!

If $K(\mathbf{x}, \mathbf{x}')$ is an inner product in **some** space \mathcal{Z} , we are good.

Example: $K(\mathbf{x}, \mathbf{x}') = e^{-\gamma \|\mathbf{x}-\mathbf{x}'\|^2}$

Infinite-dimensional \mathcal{Z} : take simple case (one dimension).

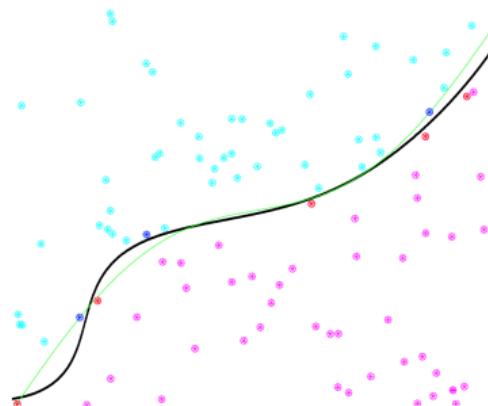
$$\begin{aligned} K(\mathbf{x}, \mathbf{x}') &= e^{-(x-x')^2} \\ &= e^{-x^2} e^{-x'^2} \underbrace{\sum_{k=0}^{\infty} \frac{2^k x^k x'^k}{k!}}_{e^{2xx'}} \end{aligned}$$

This Kernel in Action

Slightly non-separable case:

Transforming \mathcal{X} into ∞ —dimensional \mathcal{Z} .

Overkill? Count the support vectors.



Kernel Formulation SVM

Remember quadratic programming? The only difference now is:

$$\underbrace{\begin{bmatrix} y_1y_1 K(\mathbf{x}_1, \mathbf{x}_1) & y_1y_2 K(\mathbf{x}_1, \mathbf{x}_2) & \dots & y_1y_N K(\mathbf{x}_1, \mathbf{x}_N) \\ y_2y_1 K(\mathbf{x}_2, \mathbf{x}_1) & y_2y_2 K(\mathbf{x}_2, \mathbf{x}_2) & \dots & y_2y_N K(\mathbf{x}_2, \mathbf{x}_N) \\ \dots & \dots & \dots & \dots \\ y_Ny_1 K(\mathbf{x}_N, \mathbf{x}_1) & y_Ny_2 K(\mathbf{x}_N, \mathbf{x}_2) & \dots & y_Ny_N K(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}}_{\text{quadratic coefficients}}$$

Everything else is the same

The Final Hypothesis

Express $g(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{z} + b)$ in terms of $K(-, -)$.

$$\mathbf{w} = \sum_{\mathbf{z}_n \text{ in SV}} \alpha_n y_n \mathbf{z}_n \implies g(\mathbf{x}) = \text{sign} \left(\sum_{\alpha_n > 0} \alpha_n y_n K(\mathbf{x}_n, \mathbf{x}) + b \right)$$

where $b = y_m - \sum_{\alpha_n > 0} \alpha_n y_n K(\mathbf{x}_n, \mathbf{x}_m)$

for any support vector ($\alpha_m > 0$).

How Do We Know That \mathcal{Z} Exists...

... for a given $K(\mathbf{x}, \mathbf{x}')$? [valid kernel](#)

Three approaches:

1. By construction.
2. Math properties (*Mercer's condition*)
3. Who cares? ☺

Design Your Own Kernel

$K(\mathbf{x}, \mathbf{x}')$ is a valid iff:

1. It is symmetric

2. The matrix:

$$\begin{bmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & K(\mathbf{x}_1, \mathbf{x}_2) & \dots & K(\mathbf{x}_1, \mathbf{x}_N) \\ K(\mathbf{x}_2, \mathbf{x}_1) & K(\mathbf{x}_2, \mathbf{x}_2) & \dots & K(\mathbf{x}_2, \mathbf{x}_N) \\ \dots & \dots & \dots & \dots \\ K(\mathbf{x}_N, \mathbf{x}_1) & K(\mathbf{x}_N, \mathbf{x}_2) & \dots & K(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

is positive semi-definite

for any $\mathbf{x}_1, \dots, \mathbf{x}_n$ (Mercer's condition)

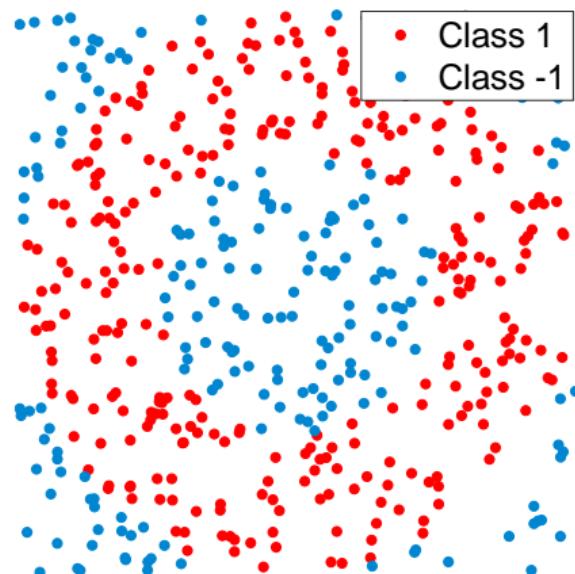
Kernel Methods: Example

- ▶ 500 observations 2 dimensions.
- ▶ Radial Basis Function kernel (RBF):

$$K(x_i, x_j) = e^{-\gamma \|x_i - x_j\|_2^2}$$

- ▶ $\gamma = 0.25$

$$\mathbf{K} = \begin{bmatrix} 1 & 0.28 & 0.52 & \dots & 0.79 \\ 0.28 & 1 & 0.03 & \dots & 0.08 \\ 0.51 & 0.03 & 1 & \dots & 0.91 \\ \dots & \dots & \dots & \dots & \dots \\ 0.79 & 0.08 & 0.91 & \dots & 1 \end{bmatrix}$$



Kernel Methods: Example

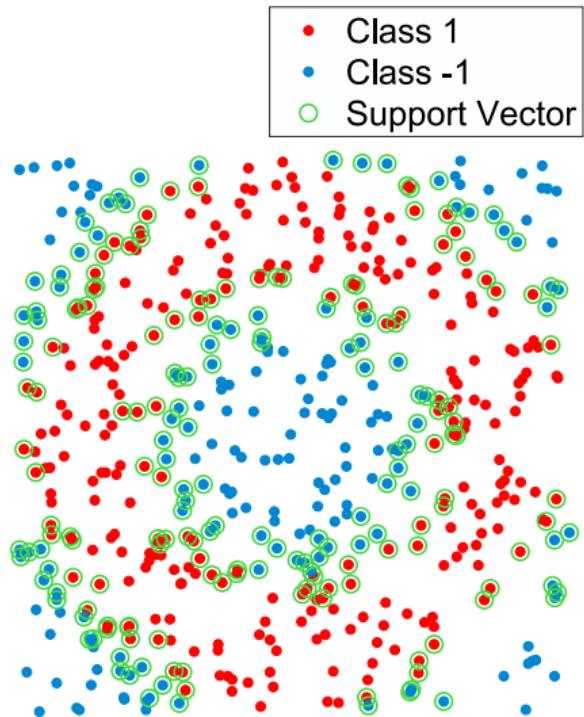
- ▶ Solve the quadratic problem (Matlab toolbox):

$$\min_{\alpha} \frac{1}{2} \alpha^T \mathbf{M} \alpha + (-\mathbf{1}^T) \alpha,$$

where

$$\mathbf{M}(i, j) = y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

$$\alpha = [0 \ 10 \ 10 \ 0 \ \dots \ 10 \ 0 \ 0]^T$$



Kernel Methods: Example

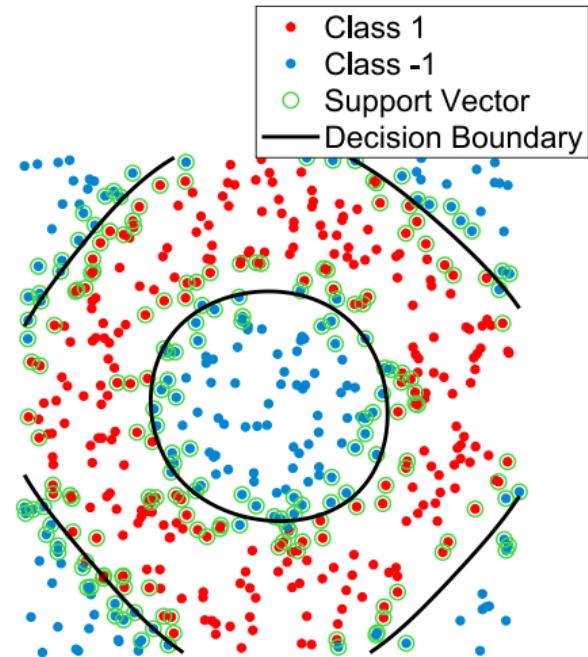
- The model:

$$b = y_m - \sum_{\alpha_n > 0} \alpha_n y_n K(\mathbf{x}_n, \mathbf{x}_m)$$

$$b = -7.78$$

$$g(\mathbf{x}) = \text{sign} \left(\sum_{\alpha_n > 0} \alpha_n y_n K(\mathbf{x}_n, \mathbf{x}) + b \right)$$

$$g(\mathbf{x}) = \text{sign} ((-1)10K(\mathbf{x}_2, \mathbf{x}) + (-1)10K(\mathbf{x}_3, \mathbf{x}) + (1)10K(\mathbf{x}_{21}, \mathbf{x}) + \dots - 7.78)$$

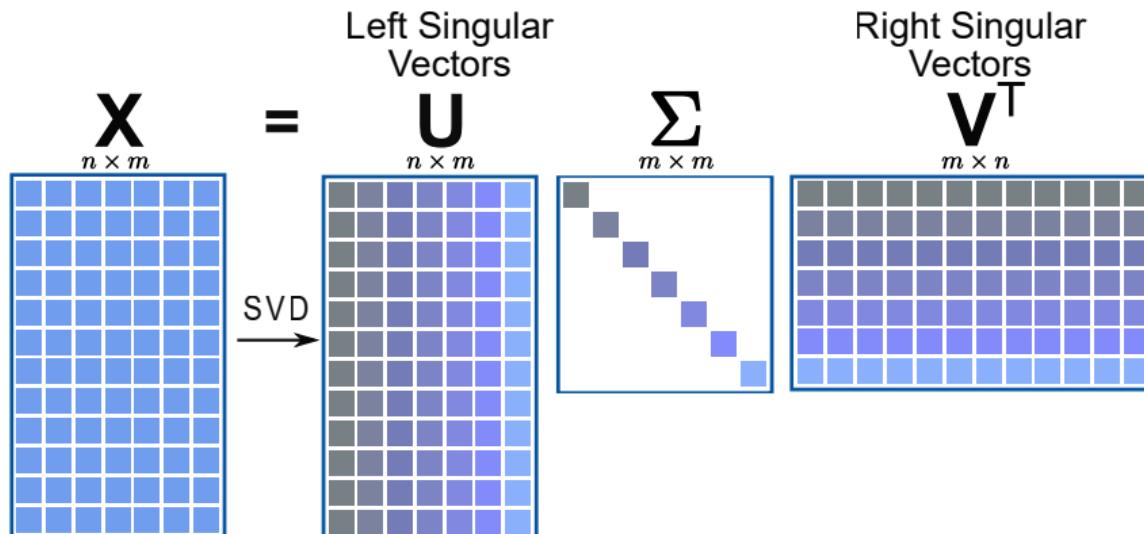


Outline

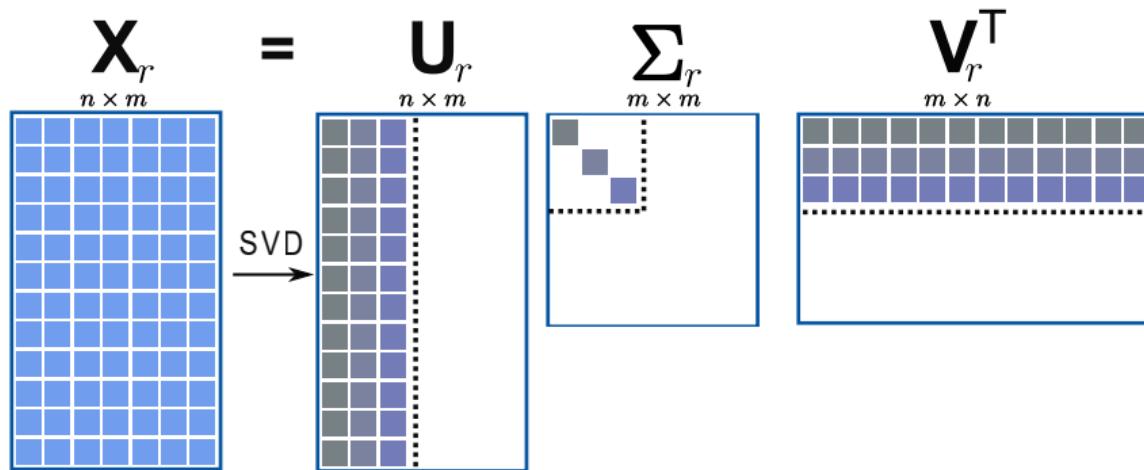
1. Aggregation.
2. Kernel Methods.
3. Randomized Singular Value Decomposition.

Randomized Singular Value Decomposition

SVD: Low-rank approximation and PCA



Randomized Singular Value Decomposition



Low-rank approximation

$$\mathbf{U}_r \mathbf{U}_r^\top \mathbf{X} = \underset{\mathbf{Z}: rank(\mathbf{Z})=r}{\operatorname{argmin}} \| \mathbf{X} - \mathbf{Z} \|_2$$

Randomized Singular Value Decomposition

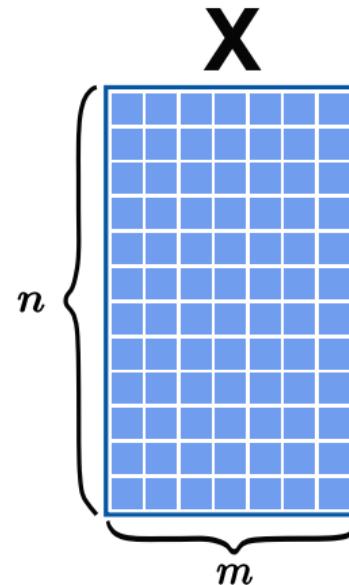
SVD → Computationally expensive $\mathcal{O}(nm^2)$

Dimensions → Increasing (4k, 8k...)

Assumption: Low intrinsic rank r .

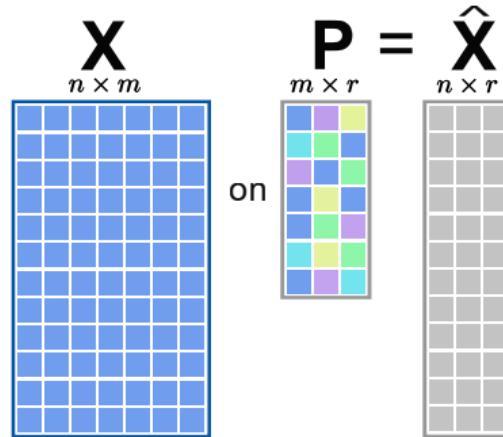
$$\mathbf{X} \approx \mathbf{U}_r \Sigma \mathbf{V}_r^\top$$

Randomly sample the column space of
 $\mathbf{X} \in \mathbb{R}^{n \times m}$.



Randomized Singular Value Decomposition

Project \mathbf{X} into $\hat{\mathbf{X}}$:



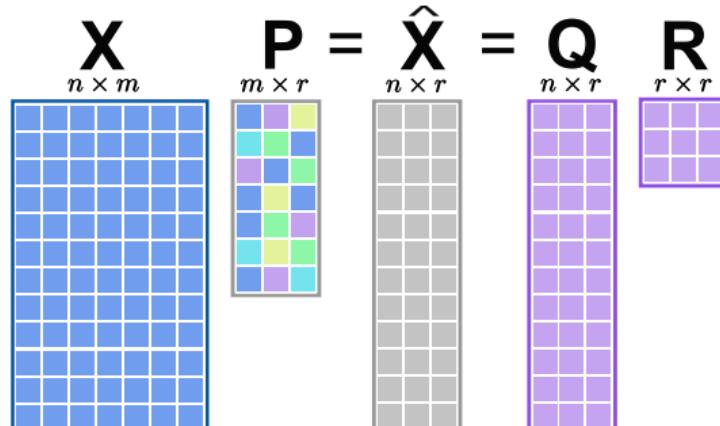
$$\begin{aligned}\| \mathbf{X} - \hat{\mathbf{U}}_r \hat{\mathbf{U}}_r^\top \mathbf{X} \|_F &\leq (1 + \epsilon) \| \mathbf{X} - \mathbf{X}_r \|_F = (1 + \epsilon) \sum_{i=r+1}^m \sigma_i^2 \\ &\leq (1 + \epsilon) \| \mathbf{X} - \mathbf{U}_r \mathbf{U}_r^\top \mathbf{X} \|_F\end{aligned}$$

Randomized Singular Value Decomposition: Step 1

- ▶ Random projection $\mathbf{P} \in \mathbb{R}^{m \times r}$.
- ▶ $\hat{\mathbf{X}} = \mathbf{XP}$, $\hat{\mathbf{X}} \in \mathbb{R}^{n \times r}$
- ▶ QR decomposition $\hat{\mathbf{X}} = \mathbf{QR}$.

\mathbf{Q} : Orthonormal basis for $\hat{\mathbf{X}}$ (and \mathbf{X}).

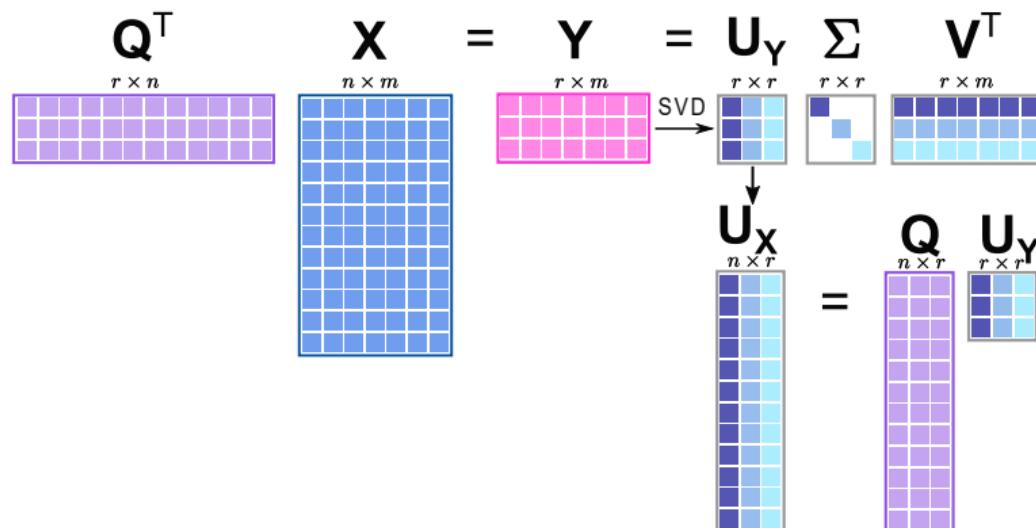
$$\mathbf{x} \approx \mathbf{QQ}^\top \mathbf{x}$$



Randomized Singular Value Decomposition: Step 2

- ▶ Project \mathbf{X} into \mathbf{Q} : $\mathbf{Y} = \mathbf{Q}^T \mathbf{X}$, $\mathbf{Y} \in \mathbb{R}^{r \times m}$.
- ▶ Compute SVD: $\mathbf{Y} = \mathbf{U}_Y \Sigma \mathbf{V}^T$
- ▶ Complexity of SVD of $\mathbf{Y} <$ Complexity of SVD of \mathbf{X}

$\Sigma \mathbf{V}^T$ Same as for \mathbf{X} and $\boxed{\mathbf{U}_X = \mathbf{Q} \mathbf{U}_Y}$



Randomized Singular Value Decomposition

Proof:

$$\begin{aligned} \mathbf{Y} &= \mathbf{Q}^T \mathbf{X} && \rightarrow \text{SVD of } \mathbf{Q} \text{ and } \mathbf{X} \\ \mathbf{Y} &= (\mathbf{U}_Q \Sigma_Q \mathbf{V}_Q^T)^T \mathbf{U}_X \Sigma_X \mathbf{V}_X^T \\ \mathbf{Y} &= \mathbf{V}_Q \Sigma_Q^T \mathbf{U}_Q^T \mathbf{U}_X \Sigma_X \mathbf{V}_X^T \end{aligned}$$

Note that $\mathbf{Q}\mathbf{Q}^T = \mathbf{I}$ (QR decomposition) $\rightarrow \mathbf{U}_Q = \mathbf{Q}$, $\mathbf{V}_Q = \mathbf{I}$ and $\Sigma_Q = \mathbf{I}$:

$$\begin{aligned} \mathbf{Y} &= \mathbf{I}^T \mathbf{Q}^T \mathbf{U}_X \Sigma_X \mathbf{V}_X^T && \rightarrow \mathbf{I} = \mathbf{I} \\ \mathbf{Y} &= \mathbf{I} \mathbf{Q}^T \mathbf{U}_X \Sigma_X \mathbf{V}_X^T && \rightarrow \mathbf{I} \mathbf{Q} = \mathbf{Q} \\ \mathbf{Y} &= \mathbf{Q}^T \mathbf{U}_X \Sigma_X \mathbf{V}_X^T && \rightarrow \text{SVD of } \mathbf{Y} \\ \mathbf{U}_Y \Sigma_Y \mathbf{V}_Y^T &= \mathbf{Q}^T \mathbf{U}_X \Sigma_X \mathbf{V}_X^T \end{aligned}$$

Then $\Sigma_Y = \Sigma_X$, $\Sigma_Y = \Sigma_X$ and $\mathbf{U}_Y = \mathbf{Q}^T \mathbf{U}_X$.

Randomized Singular Value Decomposition

Proof:

Solving for \mathbf{U}_X :

$$\begin{aligned}\mathbf{U}_Y &= \mathbf{Q}^T \mathbf{U}_X &\longrightarrow \text{Multiply both sides by } \mathbf{Q} \\ \mathbf{Q}\mathbf{U}_Y &= \mathbf{Q}\mathbf{Q}^T \mathbf{U}_X &\longrightarrow & \mathbf{Q}\mathbf{Q}^T = \mathbf{I} \\ \mathbf{Q}\mathbf{U}_Y &= \mathbf{Q}^T \mathbf{U}_X\end{aligned}$$

$$\boxed{\mathbf{U}_X = \mathbf{Q}\mathbf{U}_Y}$$

Randomized Singular Value Decomposition

Only r columns aren't always completely effective.

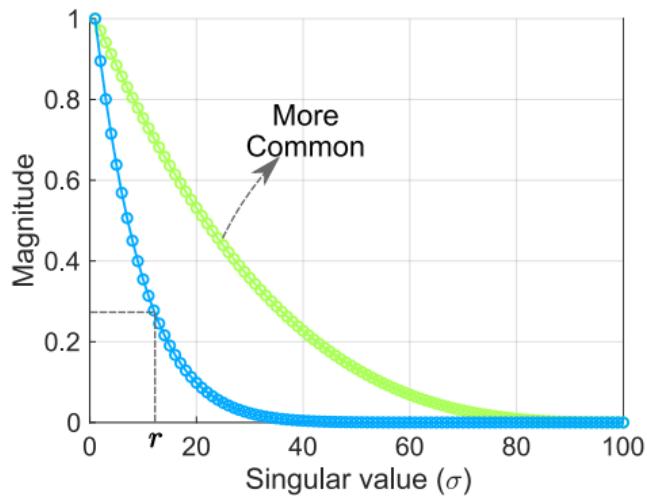
Solutions:

- ▶ Oversampling: add p columns to \mathbf{P} , $\mathbf{P} \in \mathbb{R}^{m \times (r+p)}$.
- ▶ Power iteration: $\mathbf{X}^q = (\mathbf{X}\mathbf{X}^\top)^q \mathbf{X} \longrightarrow \mathbf{X}^q = \mathbf{U}\Sigma^{2q-1}\mathbf{V}^\top$.

If \mathbf{X} is symmetric,

$$\begin{aligned}\mathbf{X} &= \mathbf{U}\Sigma\mathbf{U}^\top \\ \mathbf{X}^q &= \mathbf{U}\Sigma^q\mathbf{U}^\top\end{aligned}$$

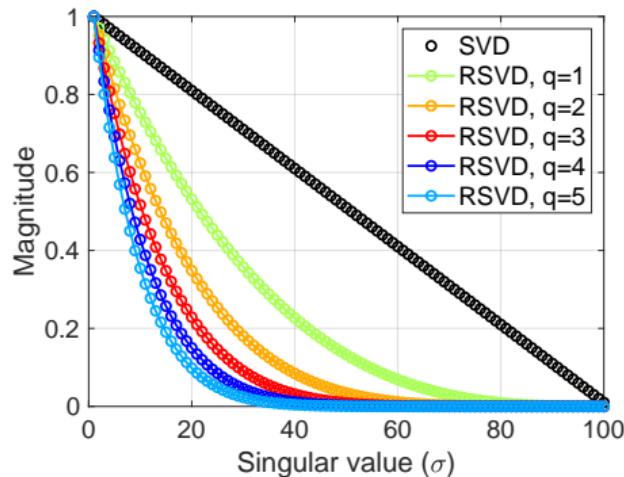
$$\|\mathbf{X}^q - \mathbf{X}_r^q\|_F = \sum_{i=r+1}^m \sigma_i^q \rightarrow \text{is much smaller}$$



Randomized Singular Value Decomposition

Power iteration: $\mathbf{X}^q = (\mathbf{X}\mathbf{X}^\top)^q \mathbf{X}$:

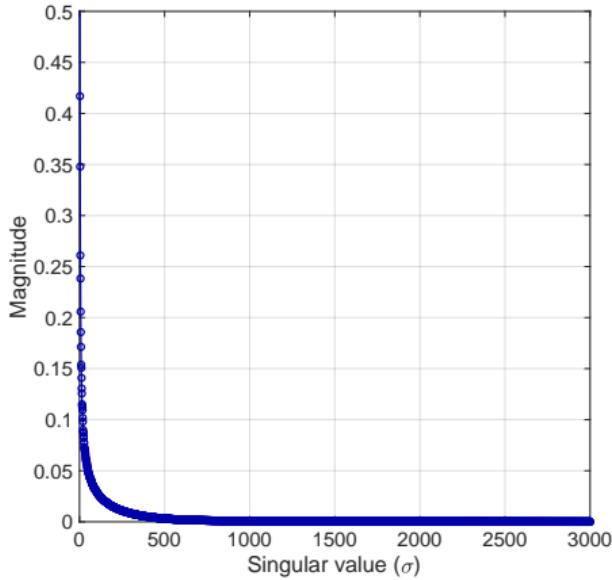
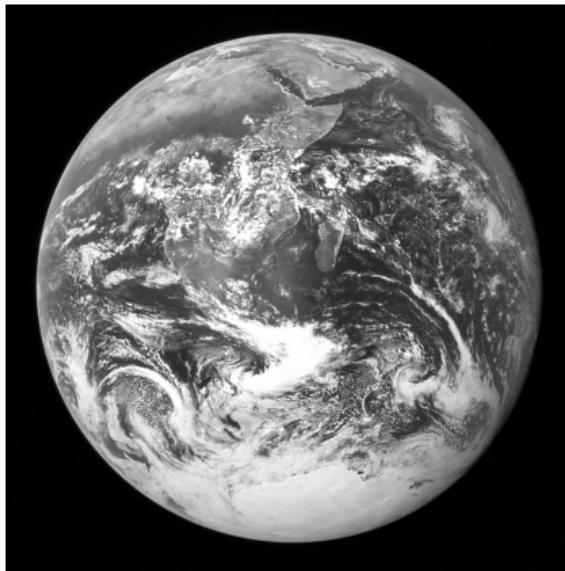
1. $\hat{\mathbf{X}} = \mathbf{X}\mathbf{P}$
2. Repeat q times $\hat{\mathbf{X}}^q = \mathbf{X}(\mathbf{X}^\top \hat{\mathbf{X}}^{q-1})$
3. QR decomposition $\hat{\mathbf{X}}^q = \mathbf{Q}\mathbf{R}$
4. Project \mathbf{X} into \mathbf{Q} : $\mathbf{Y} = \mathbf{Q}^\top \mathbf{X}$.
5. Compute SVD:
$$\mathbf{Y} = \mathbf{U}_\mathbf{Y} \Sigma \mathbf{V}^\top \longrightarrow \mathbf{U}_\mathbf{X} = \mathbf{Q} \mathbf{U}_\mathbf{Y}$$
.



- ▶ Random Matrix $\mathbf{X} \in \mathbb{R}^{1000 \times 100}$.
- ▶ Singular value spectrum of \mathbf{X}^q decays more rapidly.

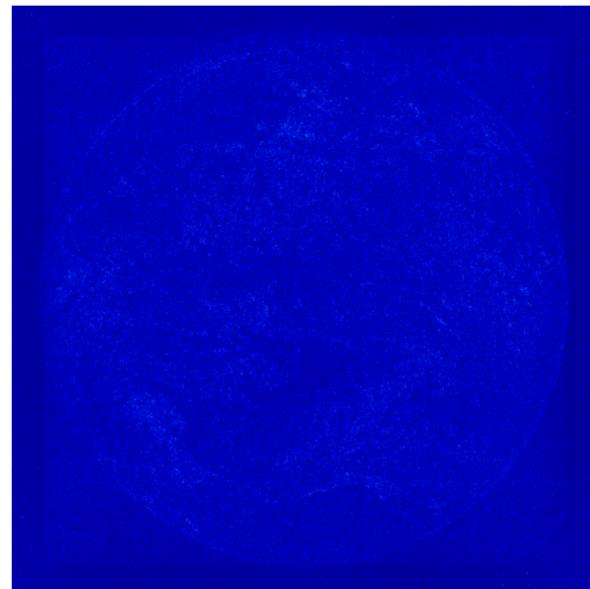
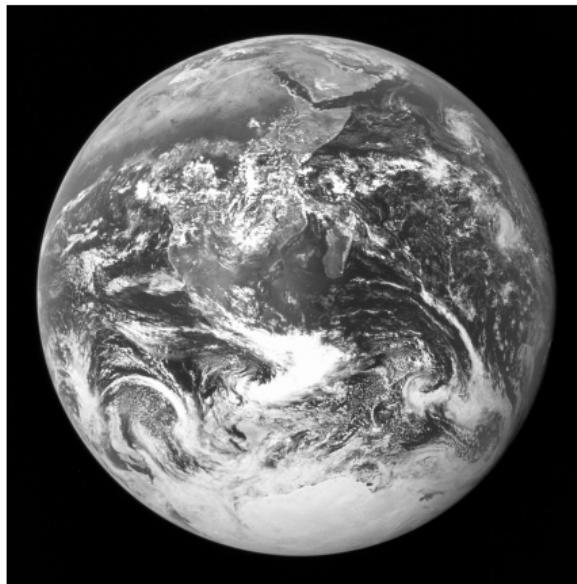
Randomized Singular Value Decomposition: Example

- ▶ Planet earth photography: 3000×3000 pixels.
- ▶ Rank: $r \approx 400$



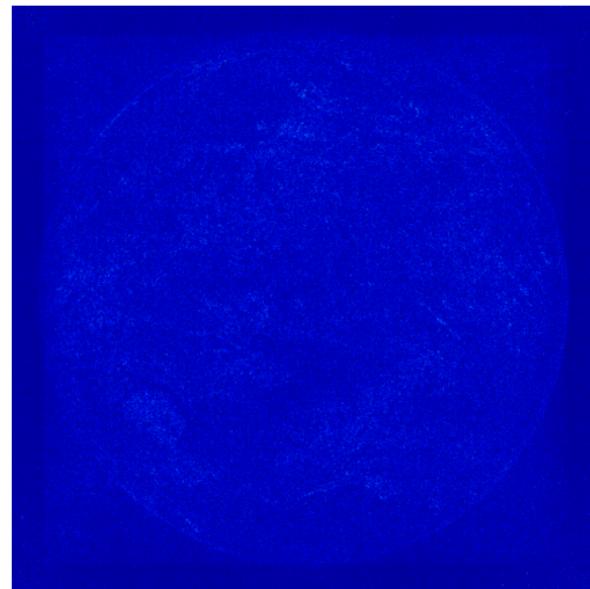
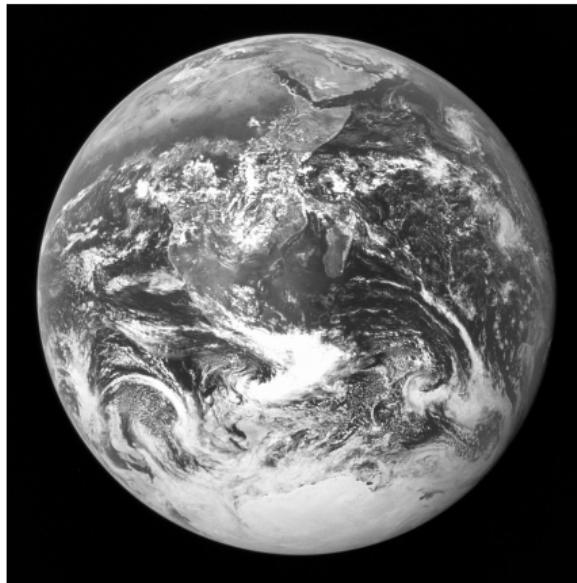
Randomized Singular Value Decomposition: Example

- ▶ SVD 400-rank approximation.
- ▶ MSE=5.7



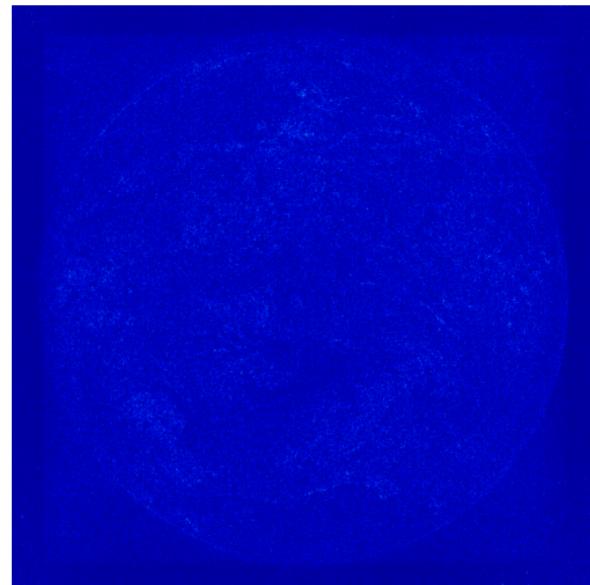
Randomized Singular Value Decomposition: Example

- ▶ Randomized SVD 400-rank approximation $\mathbf{P} \in \mathbb{R}^{3000 \times 400}$
- ▶ MSE=6.7



Randomized Singular Value Decomposition: Example

- ▶ Oversampling randomized SVD 400-rank approximation $p = 20$.
- ▶ $\mathbf{P} \in \mathbb{R}^{3000 \times (400+20)}$.
- ▶ MSE=6.3



Randomized Singular Value Decomposition: Example

- ▶ Power iterations randomized SVD 400-rank approximation $q = 2$.
- ▶ MSE=6.05

