

Introduction to Machine Learning

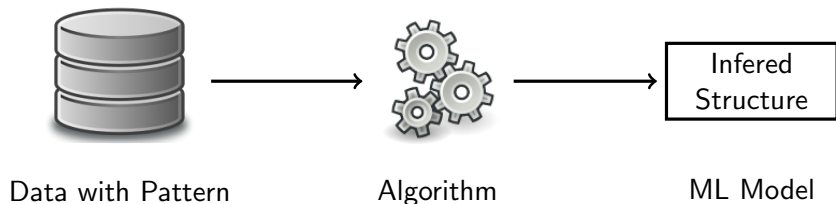
Felix Brockherde^{1,2} Kristof Schütt¹

¹Technische Universität Berlin

²Max Planck Institute of Microstructure Physics

IPAM Tutorial 2013

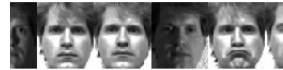
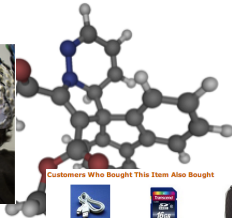
What is Machine Learning?



ML is about learning structure from data

Examples

Drug discovery



Face recognition

Google machine learning

Web Images Maps Shopping

About 129,000,000 results (0.10 second)

[Machine learning - Wikipedia, the](#)

[en.wikipedia.org/wiki/Machine_learning](#)

Machine learning, a branch of artificial intelligence, concerns the construction and study of systems that can learn from data. For example, a machine learning ...
Artificial intelligence - Supervised learning - List of machine learning ... - Wiki

[News for machine learning](#)

[DataBox API Enables Third Party Apps with Machine Learning Capabilities](#)

[ProgrammableWeb \[Blog\] - 1 day ago](#)

DataBox provides a machine learning platform that specializes in natural language processing. The DataBox API enables developers to ...

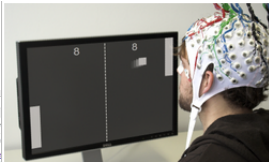
[Machine Learning Courses](#)

<https://www.coursera.org/courses>

Machine learning is the science of getting computers to act without being programmed. In the past decade, machine learning has given us self-7,847 people +112 bits

[Machine Learning - MIT OpenCourseWare](#)

ocw.mit.edu / Courses / Electrical Engineering and Computer Science 6.867 is an introductory course on machine learning which gives an overview of machine learning and advances in machine learning. has



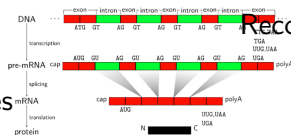
BCI

Customers Who Bought This Item Also Bought

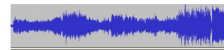
Canon PowerShot A2505 Digital Camera Compatible USB 2.0 Cable Core ... ★★★★☆ (5) \$1.99	Transcend 16GB Class 10 SDHC Flash Memory Card (TS16GSDHC10C) ★★★★★ (5,249) \$12.95	AmazonBasics Compact Camera Case (Black) ★★★★★ (67) \$4.50	SquareTrade 2-Year Protection Plan (\$75-100) ★★★★★ (505) \$12.82
---	---	--	---

Recommender systems

Search engines



DNA splice site detection



Speech recognition

Part 1: Learning Theory and Supervised ML

- ▶ Basic Ideas of Learning Theory
- ▶ Support Vector Machines
- ▶ Kernels
- ▶ Kernel Ridge Regression

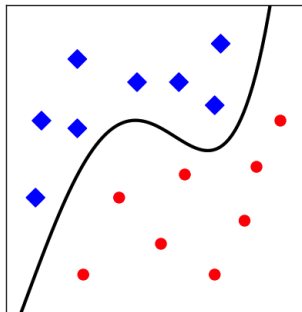
Part 2: Unsupervised ML and Application

- ▶ PCA
- ▶ Model Selection
- ▶ Feature Representation

Not covered

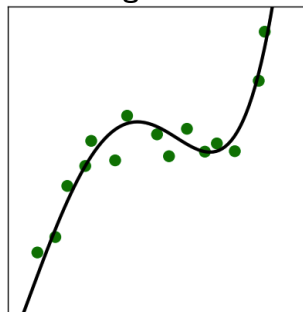
- ▶ Probabilistic Models
- ▶ Neural Networks
- ▶ Online Learning
- ▶ Reinforcement Learning
- ▶ Semi-supervised Learning
- ▶ etc.

Classification



$$y_i \in \{-1, +1\}$$

Regression



$$y_i \in \mathbb{R}$$

- ▶ Given: *Points* $X = (\mathbf{x}_1, \dots, \mathbf{x}_N)$ with $\mathbf{x}_i \in \mathbb{R}^d$ and *Labels* $Y = (y_1, \dots, y_n)$ generated by some joint probability distribution.
- ▶ Learn underlying unknown mapping $f(\mathbf{x}) = y$
- ▶ Important: Performance on unseen data

Risk minimization (RM)

Learn a model function f from examples

$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \in \mathbb{R}^d \times \mathbb{R}$ or $\{+1, -1\}$, generated from $P(\mathbf{x}, y)$

such that the expected number of errors on test data (drawn from $P(\mathbf{x}, y)$),

$$R[f] = \int \frac{1}{2} |f(\mathbf{x}) - y|^2 dP(\mathbf{x}, y),$$

is minimal.

Problem: Distribution $P(\mathbf{x}, y)$ is unknown

Empirical Risk Minimization (ERM)

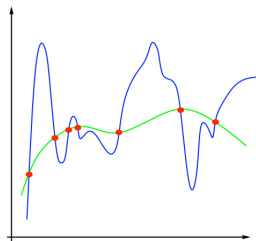
Replace the average over $P(\mathbf{x}, y)$ by average of training samples (i.e. *minimize the training error*):

$$R_{\text{emp}}[f] = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} |f(\mathbf{x}_i) - y_i|^2$$

Law of large numbers: $R_{\text{emp}}[f] \rightarrow R[f]$ as $N \rightarrow \infty$.

Question: Does $\min_f R_{\text{emp}}[f]$ give us $\min_f R[f]$ for sufficiently large N ?

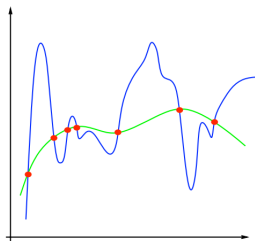
No: uniform convergence needed



Law of large numbers: $R_{\text{emp}}[f] \rightarrow R[f]$ as $N \rightarrow \infty$.

Question: Does $\min_f R_{\text{emp}}[f]$ give us $\min_f R[f]$ for sufficiently large N ?

No: uniform convergence needed



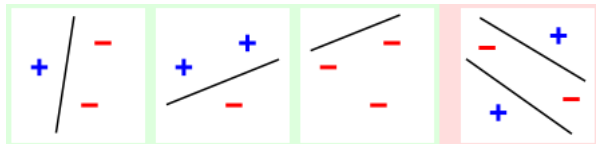
Error bound for classification

With probability of at least $1 - \eta$:

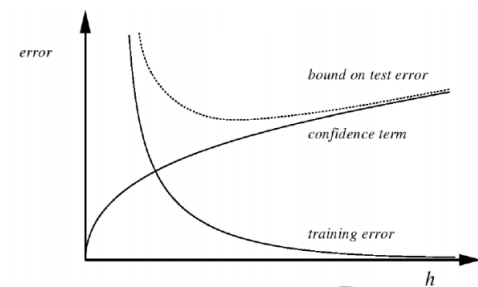
$$R[f] \leq R_{\text{emp}}[f] + \sqrt{\frac{D(\log \frac{2N}{D} + 1) - \log(\frac{\eta}{4})}{N}}$$

where D is the VC dimension (Vapnik and Chervonenkis (1971)).

Introduce structure on set of possible functions and use Structural Risk Minimization (SRM).

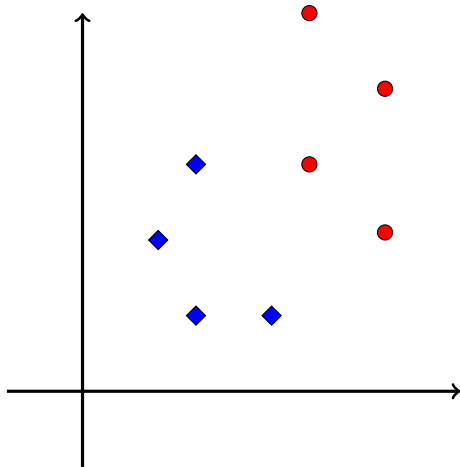


The linear function class has VC-dimension $D = 3$

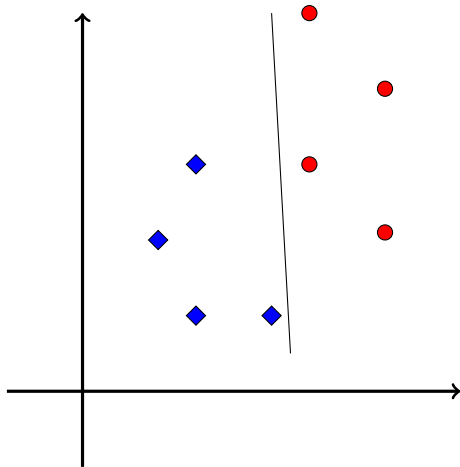


$$\min_f R_{\text{emp}}[f] + \text{Complexity}[f]$$

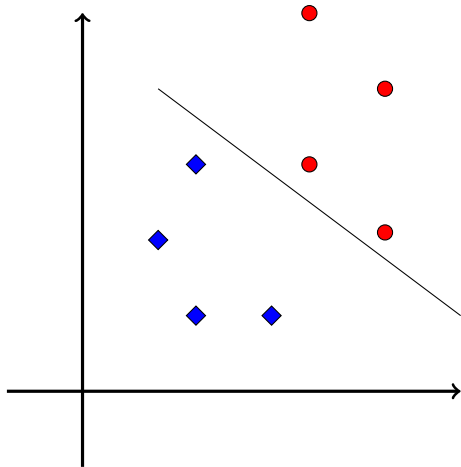
Support Vector Machines (SVM)



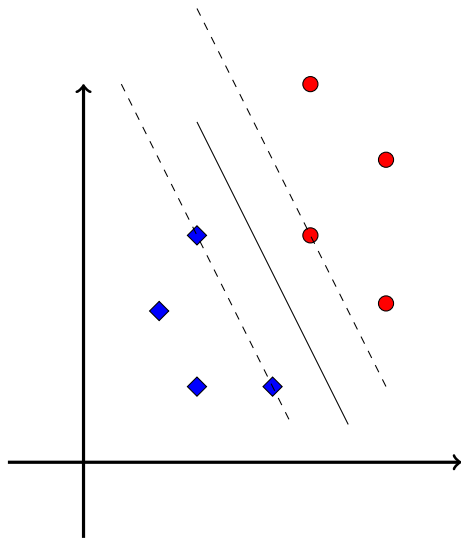
Support Vector Machines (SVM)



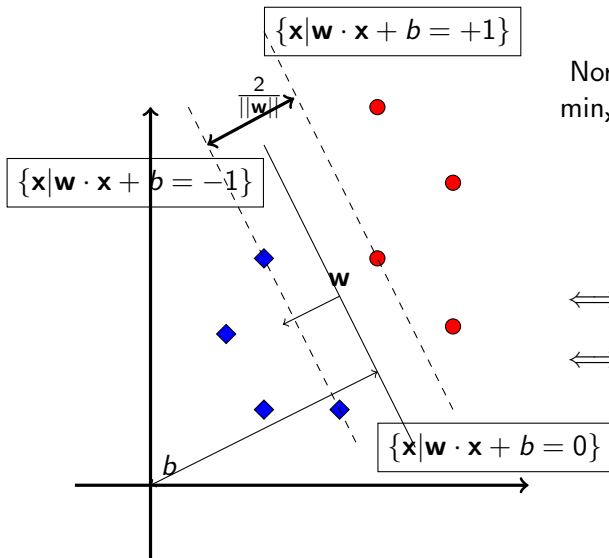
Support Vector Machines (SVM)



Support Vector Machines (SVM)



Support Vector Machines (SVM)



Normalize \mathbf{w} so that
 $\min_{x_i} \mathbf{w} \cdot \mathbf{x}_i + b = 1.$

$$\mathbf{w} \cdot \mathbf{x}_1 + b = +1$$

$$\mathbf{w} \cdot \mathbf{x}_2 + b = -1$$

$$\Leftrightarrow \mathbf{w} \cdot (\mathbf{x}_1 - \mathbf{x}_2) = 2$$

$$\Leftrightarrow \frac{\mathbf{w}}{\|\mathbf{w}\|} \cdot (\mathbf{x}_1 - \mathbf{x}_2) = \frac{2}{\|\mathbf{w}\|}$$

VC Dimension of Hyperplane Classifiers

Theorem (Cortes and Vapnik (1995))

Hyperplanes in canonical form have VC Dimension

$$D \leq \min\{R^2\|\mathbf{w}\|^2 + 1, N + 1\}$$

where R the radius of the smallest sphere containing the data.

SRM Bound:

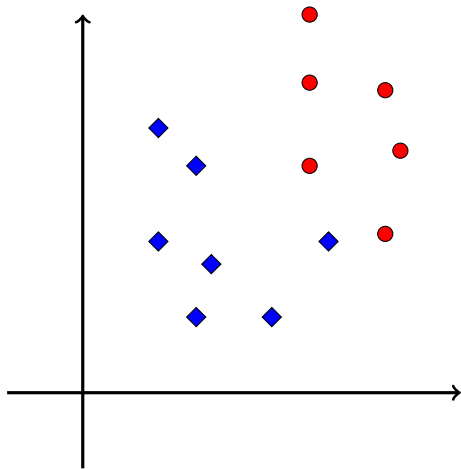
$$R[f] \leq R_{\text{emp}}[f] + \sqrt{\frac{D(\log \frac{2N}{D} + 1) - \log(\frac{\eta}{4})}{N}}$$

maximal margin = minimum $\|\mathbf{w}\|^2 \rightarrow$ good generalization, i.e. low risk:

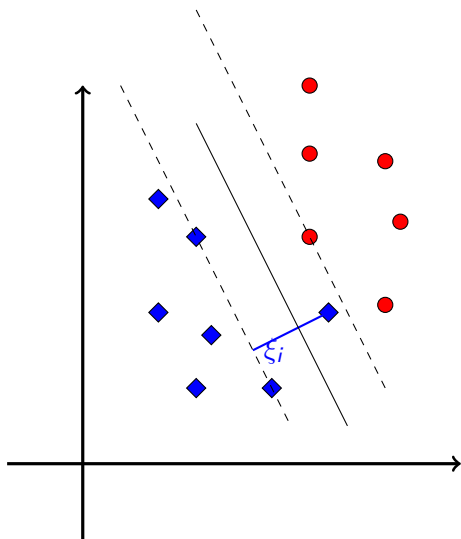
$$\min_{\mathbf{w}, b} \|\mathbf{w}\|^2$$

$$\text{subject to } y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 \quad \text{for } i = 1 \dots N$$

Slack variables



Slack variables

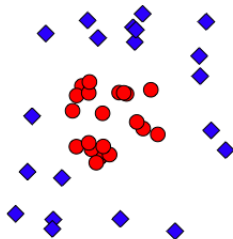


Introduce slack variables ξ_i :

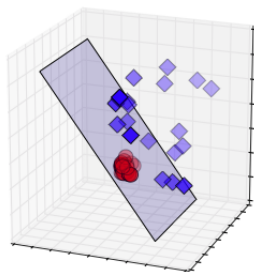
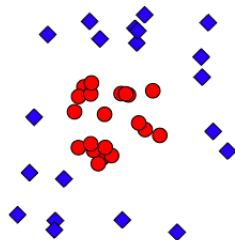
$$\min_{\mathbf{w}, b, \xi_i} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i$$

$$\text{subject to } y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$$
$$\xi_i \geq 0$$

Non-linear hyperplanes



Non-linear hyperplanes



Map into a higher dimensional feature space:

$$\Phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$(x_1, x_2) \mapsto (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

Primal

$$\min_{\mathbf{w}, b, \xi_i} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i$$

subject to $y_i(\mathbf{w} \cdot \Phi(\mathbf{x}_i) + b) \geq 1 - \xi_i$ and $\xi_i \geq 0$ for $i = 1 \dots N$

Dual

$$\max_{\alpha} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j (\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j))$$

subject to $\sum_{i=1}^N \alpha_i y_i = 0$ and $C \geq \alpha_i \geq 0$ for $i = 1 \dots N$

Data points x_i only appear in scalar products $(\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j))$.

The Kernel Trick

Replace scalar products with kernel function (Müller et al. (2001)):

$$k(x, y) = \Phi(x) \cdot \Phi(y)$$

- ▶ Compute kernel matrix $K_{ij} = k(x_i, x_j)$, **i.e. never use Φ directly**
- ▶ Underlying mapping Φ can be unknown
- ▶ Kernels can be adopted to specific task, e.g. using prior knowledge (kernels for graphs, trees, strings, ...)

Common kernels

Gaussian Kernel: $k(x, y) = \exp\left(-\frac{\|x-y\|^2}{2\sigma^2}\right)$

Linear Kernel: $k(x, y) = x \cdot y$

Polynomial Kernel: $k(x, y) = (x \cdot y + c)^d$

The Support Vectors in SVM

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j (\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)) \\ \text{subject to} \quad & \sum_{i=1}^N \alpha_i y_i = 0 \text{ and } C \geq \alpha_i \geq 0 \text{ for } i = 1 \dots N \end{aligned}$$

KKT conditions

$$y_i [\mathbf{w}\Phi(\mathbf{x}_i) + b] > 1 \implies a_i = 0 \longrightarrow x_i \text{ irrelevant}$$

$$y_i [\mathbf{w}\Phi(\mathbf{x}_i) + b] = 1 \implies \text{on/in margin} \longrightarrow x_i \text{ Support Vector}$$

Old model $f(x) = w \cdot \Phi(x_i) + b$ becomes via $w = \sum_{i=1}^N \alpha_i y_i \Phi(x_i)$:

$$f(x) = \sum_{i=1}^N \alpha_i y_i k(x_i, x) + b \longrightarrow f(x) = \sum_{x_i \in \text{SV}} \alpha_i y_i k(x_i, x) + b$$

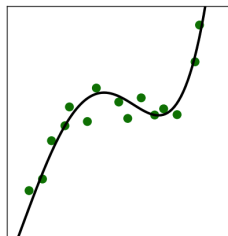
Kernel Ridge Regression (KRR)

Ridge Regression

$$\min_w \sum_{i=1}^N \|y_i - w \cdot \mathbf{x}_i\|^2 + \lambda \|w\|^2$$

Setting derivative to zero gives

$$w = \left(\lambda I + \sum_{i=1}^N x_i x_i^\top \right)^{-1} \sum_{i=1}^N y_i x_i$$



Linear Model: $f(x) = w \cdot x$

Kernelizing Ridge Regression

Setting $X = (x_1, \dots, x_N) \in \mathbb{R}^{d \times N}$ and $Y = (y_1, \dots, y_n)^T \in \mathbb{R}^N$:

$$w = (\lambda I + XX^T)^{-1} XY$$

Apply Woodbury Matrix identity:

$$w = X (X^T X + \lambda I)^{-1} Y$$

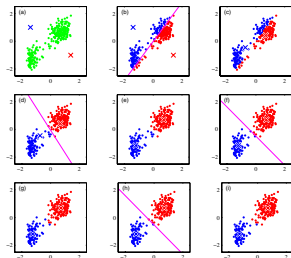
Introduce α :

$$\alpha = (K + \lambda I)^{-1} Y \quad \text{and} \quad w = \sum_{i=1}^N \Phi(\mathbf{x}_i) \alpha_i$$

$$\text{Kernel Model: } f(x) = w \cdot \Phi(x) = \sum_{i=1}^N \alpha_i k(\mathbf{x}_i, \mathbf{x})$$

Unsupervised Learning

- ▶ Learn structure from unlabeled data
- ▶ Fit an assumed model / distribution to the data
- ▶ Examples
 - ▶ clustering
 - ▶ blind source separation
 - ▶ outlier detection
 - ▶ dimensionality reduction



Principal Component Analysis (PCA)

Given centered data matrix $X = (x_1, \dots, x_N)^T \in \mathbb{R}^{N \times D}$

- ▶ best linear approximation

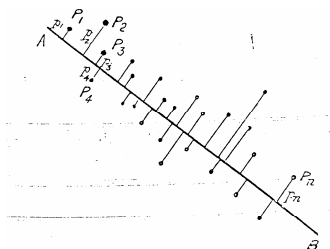
$$w_1 = \arg \min_{\|w\|=1} \{ \|X - Xww^T\|^2 \}$$

- ▶ direction of largest variance

$$w_1 = \arg \max_{\|w\|=1} \{ \|Xw\|^2 \}$$

- ▶ matrix reduction for further components

$$X_{k+1} = X_k - X_k w w^T$$



Pearson (1901)

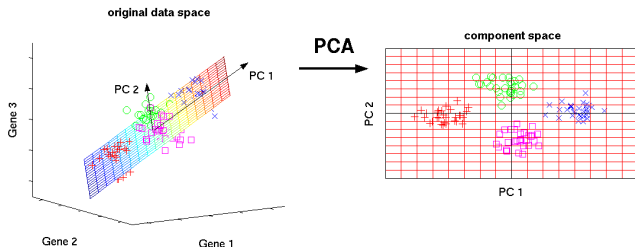
Principal Component Analysis (PCA)

Given centered data matrix $X \in \mathbb{R}^{N \times D}$, decompose correlated data matrix into uncorrelated, orthogonal *PCs*

- ▶ diagonalize covariance matrix $\Sigma = \frac{1}{N} X^T X$

$$\Sigma \mathbf{w}_k = \sigma_k^2 \mathbf{w}_k$$

- ▶ order principal components \mathbf{w}_k by variance σ_k^2
- ▶ project data to first n principal components



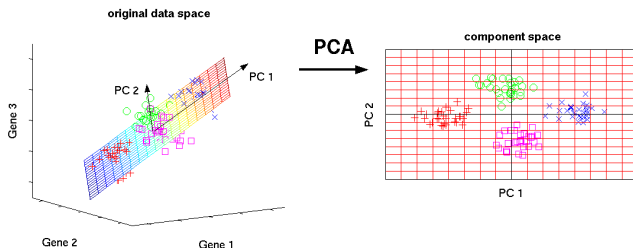
Principal Component Analysis (PCA)

Given centered data matrix $X \in \mathbb{R}^{N \times D}$, decompose correlated data matrix into uncorrelated, orthogonal *PCs*

- ▶ diagonalize covariance matrix $\Sigma = \frac{1}{N} X^T X$

$$\Sigma \mathbf{w}_k = \sigma_k^2 \mathbf{w}_k$$

- ▶ order principal components \mathbf{w}_k by variance σ_k^2
- ▶ project data to first n principal components



What about nonlinear correlations?

Kernel Principal Component Analysis (kPCA)

Transformation to feature space $X \mapsto X_f$:

$$\Sigma_f = \frac{1}{N} X_f^T X_f, K = X_f X_f^T, K_{ij} = k(x_i, x_j)$$

$$\Sigma_f \mathbf{w}_k = \sigma_k^2 \mathbf{w}_k$$

Kernel Principal Component Analysis (kPCA)

Transformation to feature space $X \mapsto X_f$:

$$\Sigma_f = \frac{1}{N} X_f^T X_f, K = X_f X_f^T, K_{ij} = k(x_i, x_j)$$

$$X_f^T X_f \mathbf{w}_k = N \sigma_k^2 \mathbf{w}_k$$

Kernel Principal Component Analysis (kPCA)

Transformation to feature space $X \mapsto X_f$:

$$\Sigma_f = \frac{1}{N} X_f^T X_f, K = X_f X_f^T, K_{ij} = k(x_i, x_j)$$

$$X_f^T X_f \mathbf{w}_k = N \sigma_k^2 \mathbf{w}_k$$

$$\Downarrow \mathbf{w}_k = X_f^T \alpha_k$$

$$X_f^T X_f X_f^T \alpha_k = N \sigma_k^2 X_f^T \alpha_k$$

Kernel Principal Component Analysis (kPCA)

Transformation to feature space $X \mapsto X_f$:

$$\Sigma_f = \frac{1}{N} X_f^T X_f, K = X_f X_f^T, K_{ij} = k(x_i, x_j)$$

$$X_f^T X_f \mathbf{w}_k = N \sigma_k^2 \mathbf{w}_k$$
$$\Downarrow \mathbf{w}_k = X_f^T \boldsymbol{\alpha}_k$$

$$X_f^T X_f X_f^T \boldsymbol{\alpha}_k = N \sigma_k^2 X_f^T \boldsymbol{\alpha}_k$$
$$\Downarrow X_f \cdot$$

$$X_f X_f^T X_f X_f^T \boldsymbol{\alpha}_k = N \sigma_k^2 X_f X_f^T \boldsymbol{\alpha}_k$$

Kernel Principal Component Analysis (kPCA)

Transformation to feature space $X \mapsto X_f$:

$$\Sigma_f = \frac{1}{N} X_f^T X_f, K = X_f X_f^T, K_{ij} = k(x_i, x_j)$$

$$X_f^T X_f \mathbf{w}_k = N \sigma_k^2 \mathbf{w}_k$$

$$\Downarrow \mathbf{w}_k = X_f^T \alpha_k$$

$$X_f^T X_f X_f^T \alpha_k = N \sigma_k^2 X_f^T \alpha_k$$

$$\Downarrow X_f \cdot$$

$$K^2 \alpha_k = N \sigma_k^2 K \alpha_k$$

Kernel Principal Component Analysis (kPCA)

Transformation to feature space $X \mapsto X_f$:

$$\Sigma_f = \frac{1}{N} X_f^T X_f, K = X_f X_f^T, K_{ij} = k(x_i, x_j)$$

$$X_f^T X_f \mathbf{w}_k = N \sigma_k^2 \mathbf{w}_k$$

$$\Downarrow \mathbf{w}_k = X_f^T \alpha_k$$

$$X_f^T X_f X_f^T \alpha_k = N \sigma_k^2 X_f^T \alpha_k$$

$$\Downarrow X_f \cdot$$

$$K^2 \alpha_k = N \sigma_k^2 K \alpha_k$$

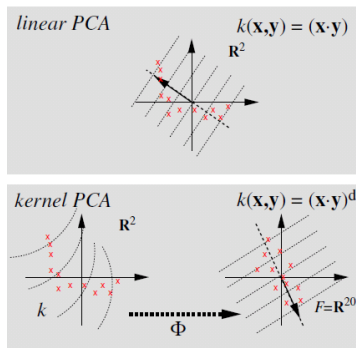
$$\Downarrow K^{-1} \cdot$$

$$K \alpha_k = N \sigma_k^2 \alpha_k$$

Kernel Principal Component Analysis (kPCA)

Projection:

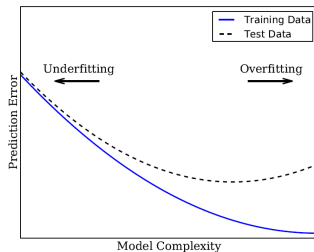
$$\begin{aligned}x_f^T w_k &= x_f^T X_f^T \alpha_k \\ &= \sum_{i=1}^N \alpha_{k,i} k(x, x_i)\end{aligned}$$



Schölkopf et al. (1997)

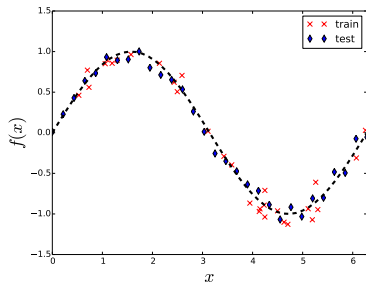
Model Selection

- ▶ Find the model that best fits the data distribution
- ▶ We can only estimate this distribution
- ▶ Consider
 - ▶ noise ratio / distribution
 - ▶ data correlation



Hyperparameters

- ▶ adjust model complexity
 - ▶ regularization, kernel parameters, etc.
- ▶ have to be tuned using examples **not** used for training
- ▶ standard solution: exhaustive search over parameter grid

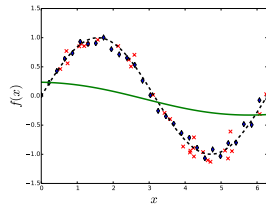
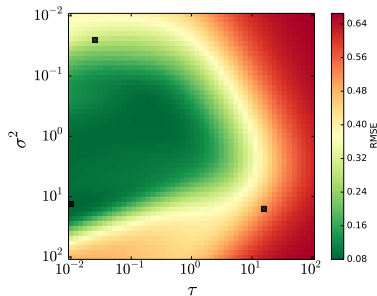
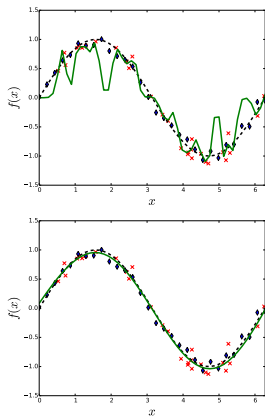


$$f(x) = \sin(x)$$

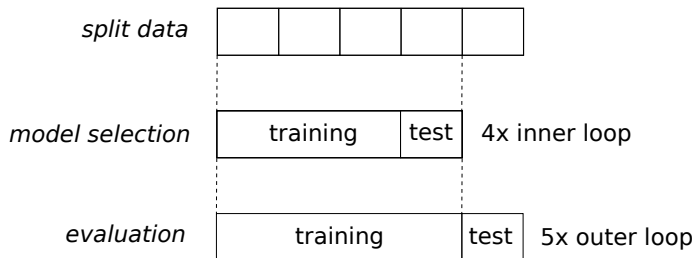
$$\bar{f}(x) = \sum_i \alpha_i \exp\left(-\frac{\|x - x_i\|^2}{\sigma^2}\right)$$

$$\alpha = (K + \tau I)^{-1} y$$

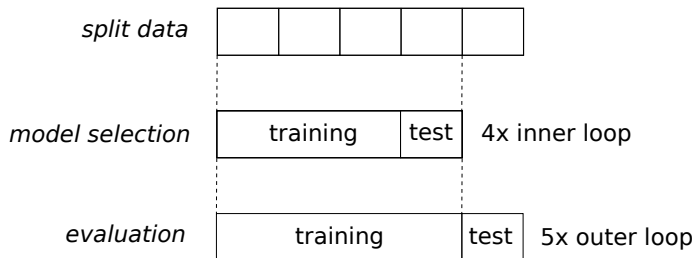
Grid Search



k-fold cross-validation



k-fold cross-validation



Don't even think about
looking at the test set!

How to represent complex objects for kernel methods?

- ▶ explicit map to vector space: $\phi : M \rightarrow \mathbb{R}^n$
 - ▶ use standard kernel (e.g., linear, polynomial, gaussian)
 $k : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ on mapped features
- ▶ direct use of kernel function: $k : M \times M \rightarrow \mathbb{R}$

Feature Representation

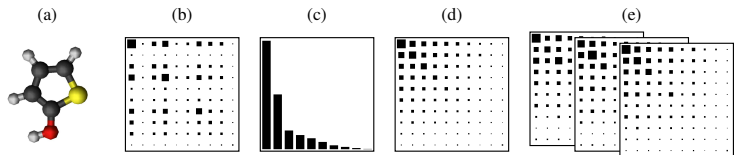
Given a physical object (molecule, crystal, etc.) and a property of interest, what is a good ML representation?

- ▶ no loss of valuable information
- ▶ support generalization
 - ▶ remove invariances
 - ▶ decompose problem
- ▶ incorporation of domain knowledge
- ▶ depends on data set, target function and learning method

Feature Representation - Molecules

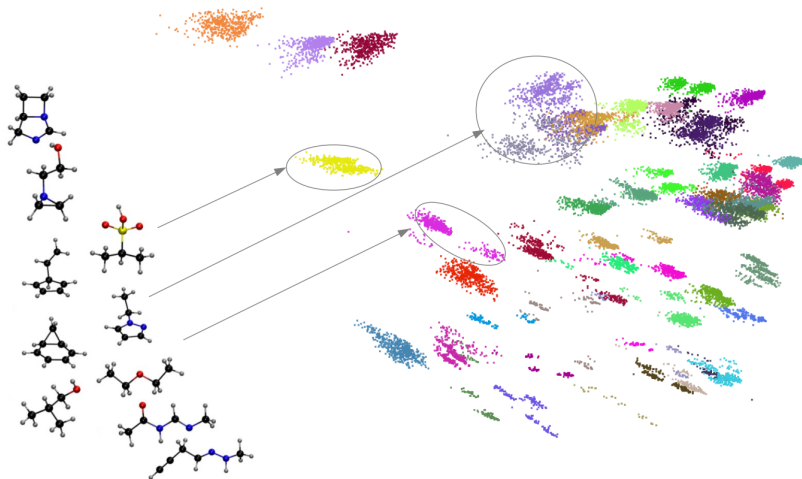
Coulomb matrix:

$$C_{ij} = \begin{cases} 0.5Z_i^{2.4} & \text{if } i = j \\ \frac{Z_i Z_j}{\|r_i - r_j\|} & \text{if } i \neq j \end{cases}$$



(Rupp et al., 2012; Montavon et al., 2012)

Feature Representation - Molecules



PCA of Coulomb matrices with atom permutations Montavon et al. (2013)

Results:

Mean predictor: **179** kcal/mol

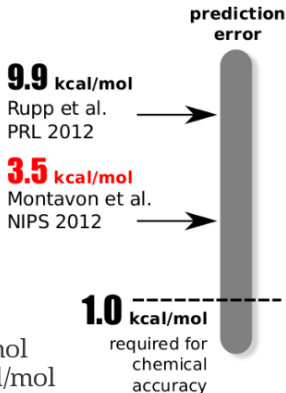
K-nearest neighbors: **70** kcal/mol

Linear regression: **20.7** kcal/mol

Deep net, backprop: **11.8** kcal/mol

Gaussian kernel ridge regression: **8.7** kcal/mol

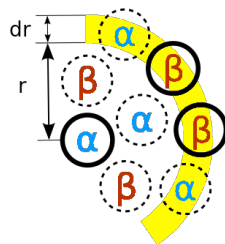
Deep net, backprop + permutations: **3.5** kcal/mol



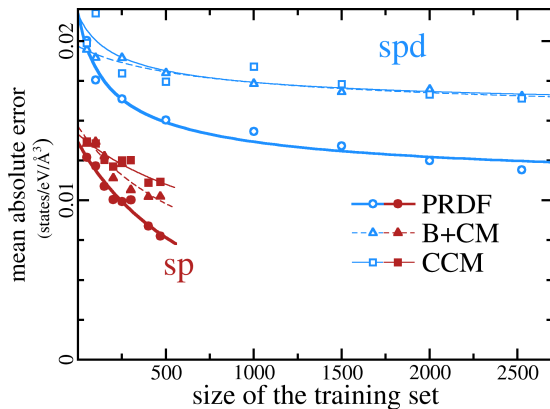
G. Montavon, K. Hansen, S. Fazli, M. Rupp, F. Biegler, A. Ziehe, A. Tkatchenko, O. A. von Lilienfeld, K.-R. Müller, [Learning Invariant Representations of Molecules for Atomization Energy Prediction](#), Advances in Neural Information Processing Systems (NIPS), 2012

Feature Representation - Crystals

element pair	r_1	\dots	r_n
$\alpha \alpha$	$g_{\alpha\alpha}(r_1)$	\dots	$g_{\alpha\alpha}(r_n)$
$\alpha \beta$	$g_{\alpha\beta}(r_1)$	\dots	$g_{\alpha\beta}(r_n)$
$\beta \alpha$	$g_{\beta\alpha}(r_1)$	\dots	$g_{\beta\alpha}(r_n)$
$\beta \beta$	$g_{\beta\beta}(r_1)$	\dots	$g_{\beta\beta}(r_n)$



Learning curve of DOS_{fermi} predictions



K.T. Schütt, H. Glawe, F. Brockherde, A. Sanna, K.-R. Müller, E.K.U. Gross, How to represent crystal structures for machine learning: towards fast prediction of electronic properties, arXiv, 2013

... has been successfully applied to various research fields.

Machine Learning ...

... has been successfully applied to various research fields.

... is based on statistical learning theory.

... has been successfully applied to various research fields.

... is based on statistical learning theory.

... provides fast and accurate predictions on previously unseen data.

- ... has been successfully applied to various research fields.
- ... is based on statistical learning theory.
- ... provides fast and accurate predictions on previously unseen data.
- ... is able to model non-linear relationships of high-dimensional data.

- ... has been successfully applied to various research fields.
- ... is based on statistical learning theory.
- ... provides fast and accurate predictions on previously unseen data.
- ... is able to model non-linear relationships of high-dimensional data.

Feature representation is key!

- Cortes, C. and Vapnik, V. (1995). Support-vector networks. *Machine learning*, 20(3):273–297.
- Montavon, G., Hansen, K., Fazli, S., Rupp, M., Biegler, F., Ziehe, A., Tkatchenko, A., Lilienfeld, A. V., and Müller, K.-R. (2012). Learning invariant representations of molecules for atomization energy prediction. In *Advances in Neural Information Processing Systems*, pages 449–457.
- Montavon, G., Rupp, M., Gobre, V., Vazquez-Mayagoitia, A., Hansen, K., Tkatchenko, A., Müller, K.-R., and von Lilienfeld, O. A. (2013). Machine learning of molecular electronic properties in chemical compound space. *arXiv preprint arXiv:1305.7074*.
- Müller, K.-R., Mika, S., Ratsch, G., Tsuda, K., and Scholkopf, B. (2001). An introduction to kernel-based learning algorithms. *Neural Networks, IEEE Transactions on*, 12(2):181–201.
- Pearson, K. (1901). Liii. on lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(11):559–572.
- Rupp, M., Tkatchenko, A., Müller, K.-R., and von Lilienfeld, O. A. (2012). Fast and accurate modeling of molecular atomization energies with machine learning. *Physical Review Letters*, 108(5):058301.
- Schölkopf, B., Smola, A., and Müller, K.-R. (1997). Kernel principal component analysis. In *Artificial Neural Networks—ICANN'97*, pages 583–588. Springer.
- Vapnik, V. N. and Chervonenkis, A. Y. (1971). On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability & Its Applications*, 16(2):264–280.