

Introduction to Machine Learning

Non-linear prediction with kernels (continued)

Prof. Andreas Krause
Learning and Adaptive Systems (las.ethz.ch)

The „Kernel Trick“

- Express problem s.t. it only depends on inner products
- Replace inner products by kernels

$$\mathbf{x}_i^T \mathbf{x}_j \quad \Rightarrow \quad k(\mathbf{x}_i, \mathbf{x}_j)$$

- This „trick“ is very widely applicable!

Definition: kernel functions

- Data space X
- A **kernel** is a function $k : X \times X \rightarrow \mathbb{R}$ satisfying
- 1) **Symmetry**: For any $\mathbf{x}, \mathbf{x}' \in X$ it must hold that

$$k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$$

- 2) **Positive semi-definiteness**: For any n , any set $S = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq X$, the kernel (Gram) matrix

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{pmatrix}$$

must be positive semi-definite

Examples of kernels on \mathbb{R}^d

- Linear kernel:

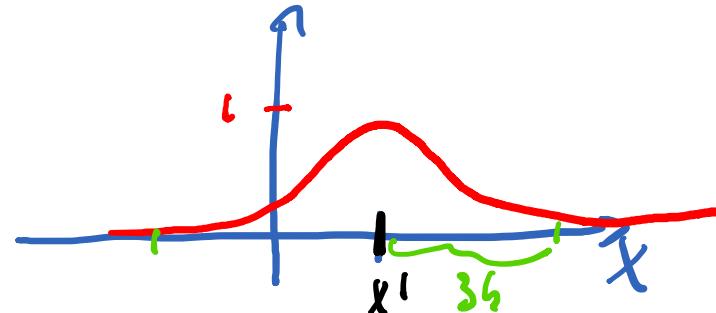
$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$$

- Polynomial kernel:

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + 1)^d$$

- Gaussian (RBF,
squared exp. kernel):

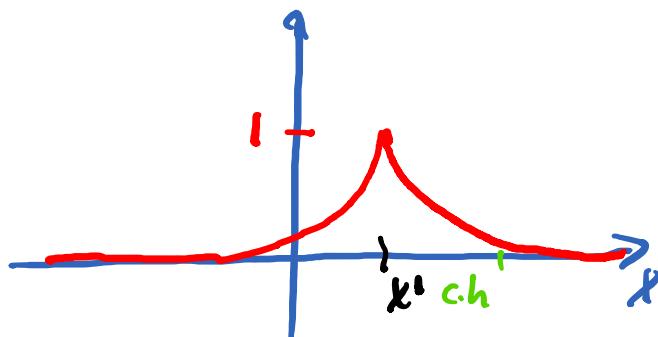
$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_2^2 / h^2)$$



Bandwidth/
Lengthscale
parameter

- Laplacian kernel:

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_1 / h)$$



Kernel engineering (composition rules)

- Suppose we have two kernels

$$k_1 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

$$k_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

defined on data space X

- Then the following functions are valid kernels:

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') k_2(\mathbf{x}, \mathbf{x}')$$

$$k(\mathbf{x}, \mathbf{x}') = c k_1(\mathbf{x}, \mathbf{x}') \text{ for } c > 0$$

$$k(\mathbf{x}, \mathbf{x}') = f(k_1(\mathbf{x}, \mathbf{x}'))$$

Sps. $\mathcal{V} : \mathcal{Z} \rightarrow \mathcal{X}$

$k(\mathbf{z}, \mathbf{z}') = k_1(\mathcal{V}(\mathbf{z}), \mathcal{V}(\mathbf{z}'))$
is valid kernel.

Pf: $k_1(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$

$k(\mathbf{z}, \mathbf{z}') = \phi(\mathcal{V}(\mathbf{z}))^T \phi(\mathcal{V}(\mathbf{z}'))$

where f is a polynomial with positive coefficients or
the exponential function

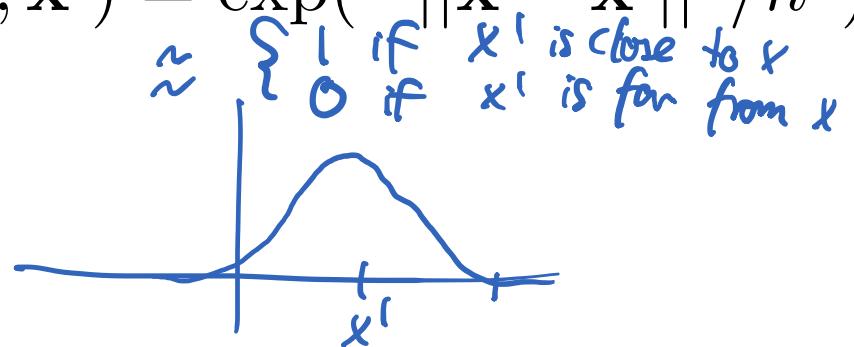
Kernels as *similarity functions*

- Recall Perceptron (and SVM) classification rule:

$$y = \text{sign} \left(\sum_{i=1}^n \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}) \right) \approx \text{sign} \left(\sum_{i=1}^n \alpha_i y_i [x \text{ "close" to } x_i] \right)$$

$\underbrace{\sum_{i=1}^n \alpha_i y_i k(\mathbf{x}_i, \mathbf{x})}_{f(x)}$

- Consider Gaussian kernel $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2/h^2)$



K-NN vs. Kernel Perceptron

- k-Nearest Neighbor:

$$y = \text{sign} \left(\sum_{i=1}^n y_i [\mathbf{x}_i \text{ among } k \text{ nearest neighbors of } \mathbf{x}] \right)$$

- Kernel Perceptron:

$$y = \text{sign} \left(\sum_{i=1}^n y_i \alpha_i k(\mathbf{x}_i, \mathbf{x}) \right)$$

Comparison: k-NN vs Kernelized Perceptron

<i>Method</i>	<i>k</i> -NN	<i>Kernelized Perceptron</i>
Advantages	No training necessary	Optimized weights can lead to improved performance Can capture „global trends“ with suitable kernels Depends on „wrongly classified“ examples only
Disadvantages	Depends on all data → inefficient	Training requires optimization

Parametric vs nonparametric learning

- Parametric models have finite set of parameters
- Example: Linear regression, linear Perceptron, ...
- Nonparametric models grow in complexity with the size of the data
 - Potentially much more expressive
 - But also more computationally complex – Why?
- Example: Kernelized Perceptron, k-NN, ...
- Kernels provide a principled way of deriving non-parametric models from parametric ones

Where are we?

- We've seen how to kernelize the perceptron
- Discussed properties of kernels, and seen examples
- Next question:
 - Can we use the kernel trick beyond the perceptron?

Kernelized SVM

- The support vector machine

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \max \{0, \underbrace{1 - y_i \mathbf{w}^T \mathbf{x}_i}_{\text{red}}\} + \underbrace{\lambda \|\mathbf{w}\|_2^2}_{\text{blue}}$$

can also be kernelized

How to kernelize the objective?

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \max \left\{ 0, 1 - y_i \mathbf{w}^T \mathbf{x}_i \right\} + \lambda \|\mathbf{w}\|_2^2$$

(*)

$\boxed{w = \sum_{j=1}^n \alpha_j y_j x_j}$

(*) = $\max \left\{ 0, 1 - y_i \left(\sum_{j=1}^n \alpha_j y_j x_j \right)^T x_i \right\}$

- $\max \left\{ 0, 1 - y_i \underbrace{\sum_{j=1}^n \alpha_j y_j}_{k(x_i, x_j)} \underbrace{(x_j^T x_i)} \right\}$

$\Rightarrow \max \left\{ 0, \alpha^T k_i \right\} \quad \text{where } k_i^T = [y_i k(x_i, x_1), \dots, y_n k(x_i, x_n)]$

How to kernelize the regularizer?

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \max\{0, 1 - y_i \mathbf{w}^T \mathbf{x}_i\} + \underbrace{\lambda \|\mathbf{w}\|_2^2}_{(\star)}$$

$$(\star) = \lambda \mathbf{w}^T \mathbf{w} = \lambda \left(\sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \right)^T \left(\sum_{j=1}^n \alpha_j y_j \mathbf{x}_j \right)$$

$$= \lambda \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \underbrace{\mathbf{x}_i^T \mathbf{x}_j}_{k(\mathbf{x}_i, \mathbf{x}_j)}$$

$$= \lambda \boldsymbol{\alpha}^T D_y K D_y \boldsymbol{\alpha}$$

$$\begin{cases} \boldsymbol{\alpha}^T = [\alpha_1, \dots, \alpha_n] \\ D_y = \begin{pmatrix} y_1 & & 0 \\ & \ddots & \\ 0 & & y_n \end{pmatrix} \end{cases}$$

$$K = \begin{pmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ \vdots & & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{pmatrix}$$

Learning & prediction with kernel classifier

- Learning: Solve the problem

Perceptron: $\arg \min_{\alpha} \frac{1}{n} \sum_{i=1}^n \max\{0, -y_i \alpha^T \mathbf{k}_i\}$ Or:

SVM: $\arg \min_{\alpha} \frac{1}{n} \sum_{i=1}^n \max\{0, 1 - y_i \alpha^T \mathbf{k}_i\} + \lambda \alpha^T \mathbf{D_y K D_y} \alpha$

$$\mathbf{k}_i = [y_1 k(\mathbf{x}_i, \mathbf{x}_1), \dots, y_n k(\mathbf{x}_i, \mathbf{x}_n)]$$

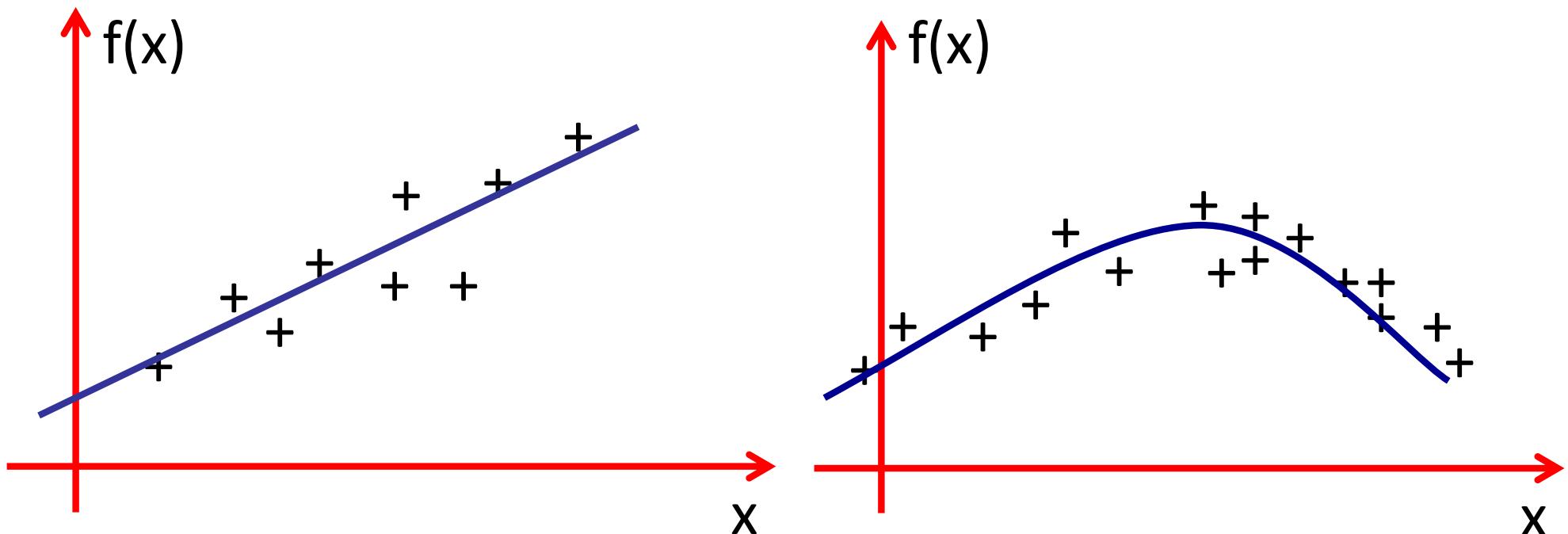
- Prediction: For data point \mathbf{x} predict label y as

$$\hat{y} = \text{sign} \left(\sum_{i=1}^n \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}) \right)$$

Demo: Kernelized SVM

Kernelized Linear Regression

- From linear to **nonlinear regression**:



- Can also kernelize linear regression
- Predictor has the form

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})$$

Example: Kernelized linear regression

- Original (**parametric**) linear optimization problem

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_2^2$$

- Similar as in perceptron, optimal $\hat{\mathbf{w}}$ lies in span of data:

$$\hat{\mathbf{w}} = \sum_{i=1}^n \alpha_i \underline{\mathbf{x}_i}$$

Kernelizing linear regression

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \underbrace{\left(\mathbf{w}^T \mathbf{x}_i - y_i \right)^2}_{(*)} + \lambda \|\mathbf{w}\|_2^2$$

$$\hat{\mathbf{w}} = \sum_{i=1}^n \alpha_i \mathbf{x}_i$$

$$(*) = \sum_{j=1}^n \alpha_j \underbrace{\mathbf{x}_j \mathbf{x}_j^T \mathbf{x}_i}_{k(\mathbf{x}_i, \mathbf{x}_j)}$$

$$k_i \in \begin{pmatrix} k(\mathbf{x}_i, \mathbf{x}_1) \\ \vdots \\ k(\mathbf{x}_i, \mathbf{x}_n) \end{pmatrix}$$

$$= \boldsymbol{\alpha}^T k_i$$

$$k = \begin{pmatrix} k_1 & \cdots & k_n \end{pmatrix}$$

$$(\text{LSS}) \equiv \hat{\boldsymbol{\alpha}} = \underset{\boldsymbol{\alpha}}{\operatorname{argmin}} \quad \frac{1}{n} \sum_{i=1}^n \underbrace{\left(\boldsymbol{\alpha}^T k_i - y_i \right)^2}_{\|\boldsymbol{\alpha}^T K - y\|_2^2} + \lambda \boldsymbol{\alpha}^T K \boldsymbol{\alpha}$$

Kernelized linear regression

$$\hat{\alpha} = \arg \min_{\alpha_{1:n}} \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^n \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) - y_i \right)^2 + \lambda \alpha^T \mathbf{K} \alpha \quad \mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{pmatrix}$$

Learning & Predicting with KLR

- Learning: Solve least squares problem

$$\hat{\alpha} = \arg \min_{\alpha} \frac{1}{n} \|\alpha^T \mathbf{K} - \mathbf{y}\|_2^2 + \lambda \alpha^T \mathbf{K} \alpha$$

Closed-form solution: $\hat{\alpha} = (\mathbf{K} + n\lambda \mathbf{I})^{-1} \mathbf{y}$

- Prediction: For data point \mathbf{x} predict response y as

$$\hat{y} = \sum_{i=1}^n \hat{\alpha}_i k(\mathbf{x}_i, \mathbf{x})$$

Demo: Kernelized linear regression

KLR for the linear kernel

- What if $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$?

$$f(x) = \sum_{i=1}^n \alpha_i k(x_i, x) = \left(\sum_{i=1}^n \alpha_i x_i^T \right) x = w^T x$$

Application: semi-parametric regression

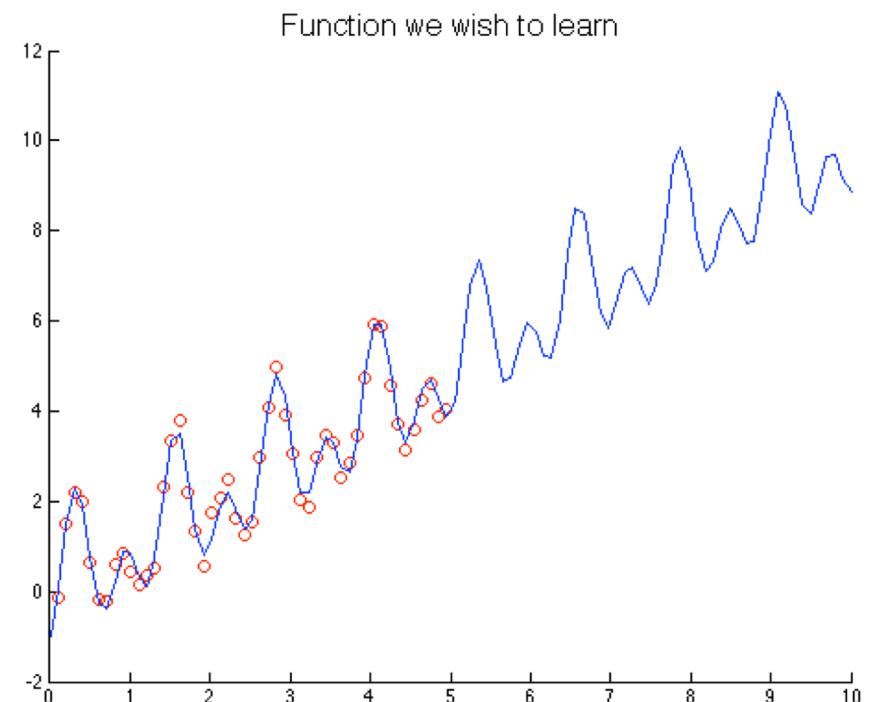
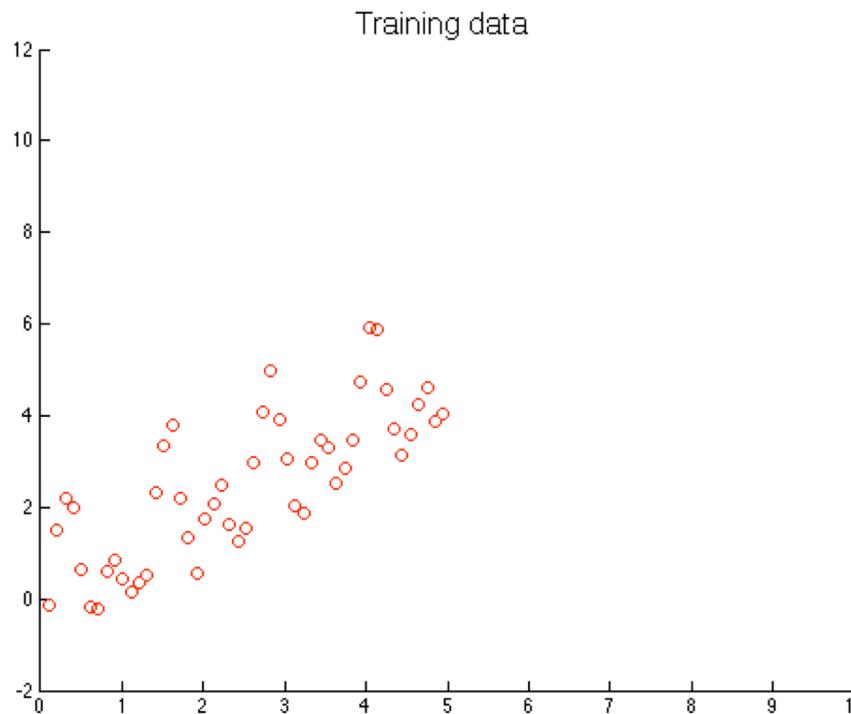
- Often, parametric models are too „rigid“, and non-parametric models fail to extrapolate
- Solution:** Use additive combination of linear and non-linear kernel function

$$k(\mathbf{x}, \mathbf{x}') = c_1 \exp(-\|\mathbf{x} - \mathbf{x}'\|_2^2/h^2) + c_2 \mathbf{x}^T \mathbf{x}'$$

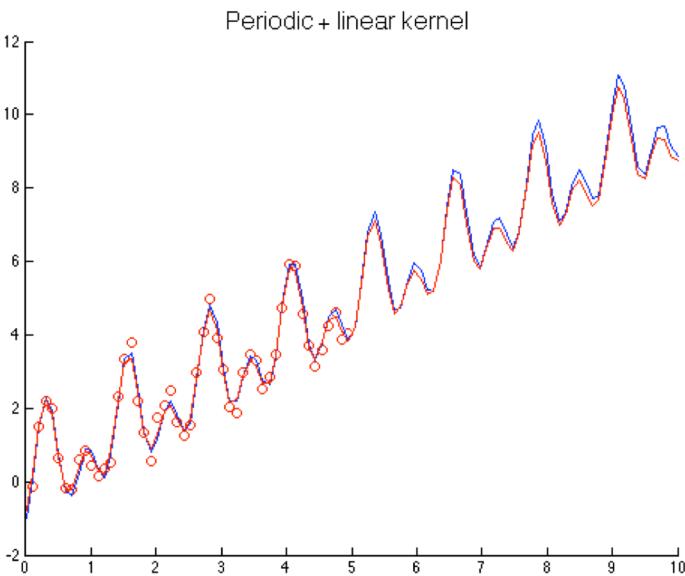
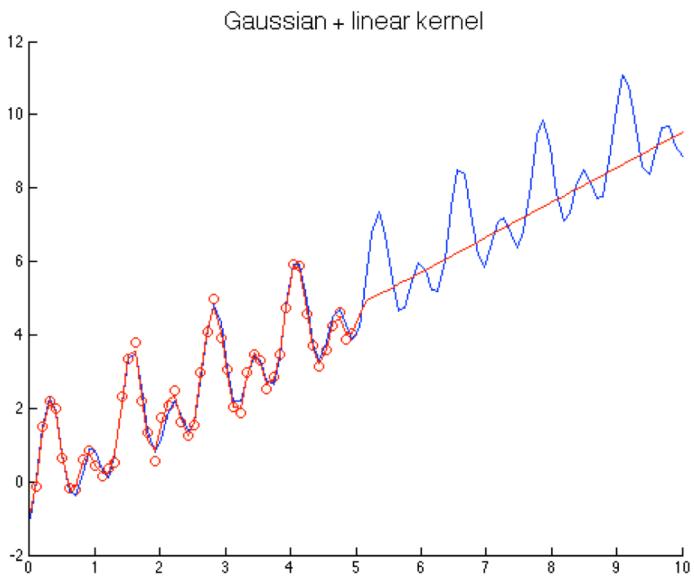
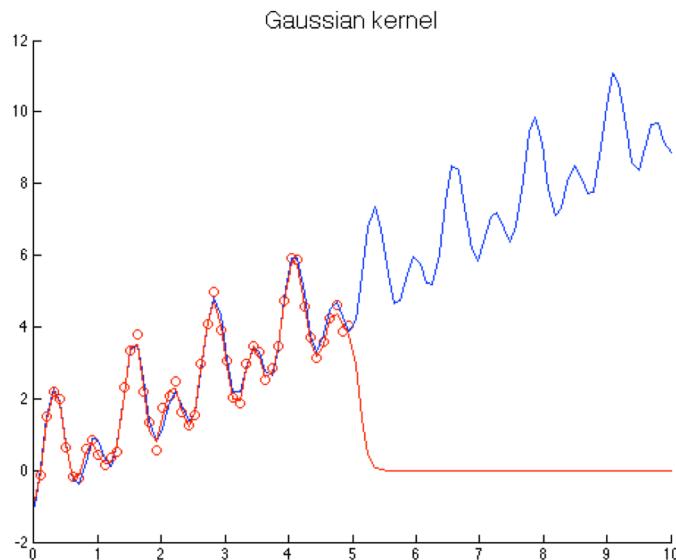
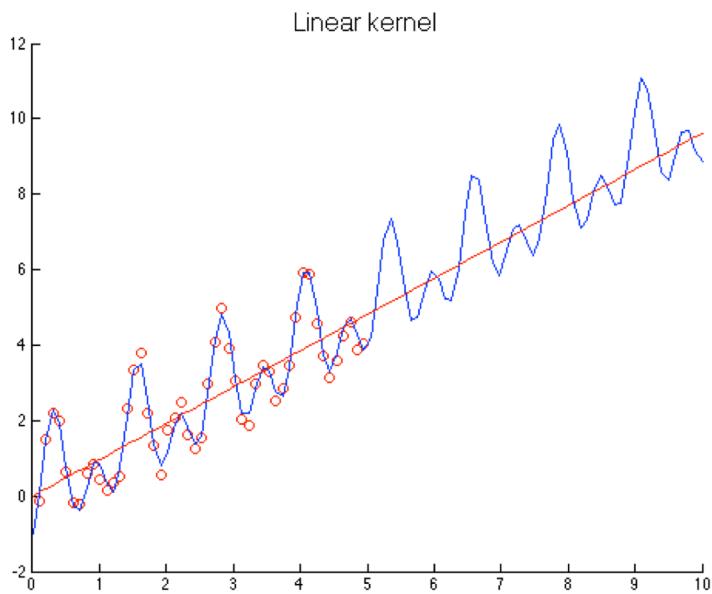
$$\begin{aligned}\hookrightarrow f(x) &= \sum_{i=1}^n \alpha_i k(x_i, x) = \sum_{i=1}^n \left(\alpha_i \left(c_1 \exp(-\|x_i - x\|_2^2/h^2) \right) + c_2 x_i^T x \right) \\ &= c_1 \sum_{i=1}^n \alpha_i \exp(-\|x_i - x\|_2^2/h^2) + \underbrace{\left(c_2 \sum_{i=1}^n \alpha_i x_i \right)^T}_{w^T} x \\ &= f_p(x) + w^T x\end{aligned}$$

Demo: Semi-parametric KLR

Example

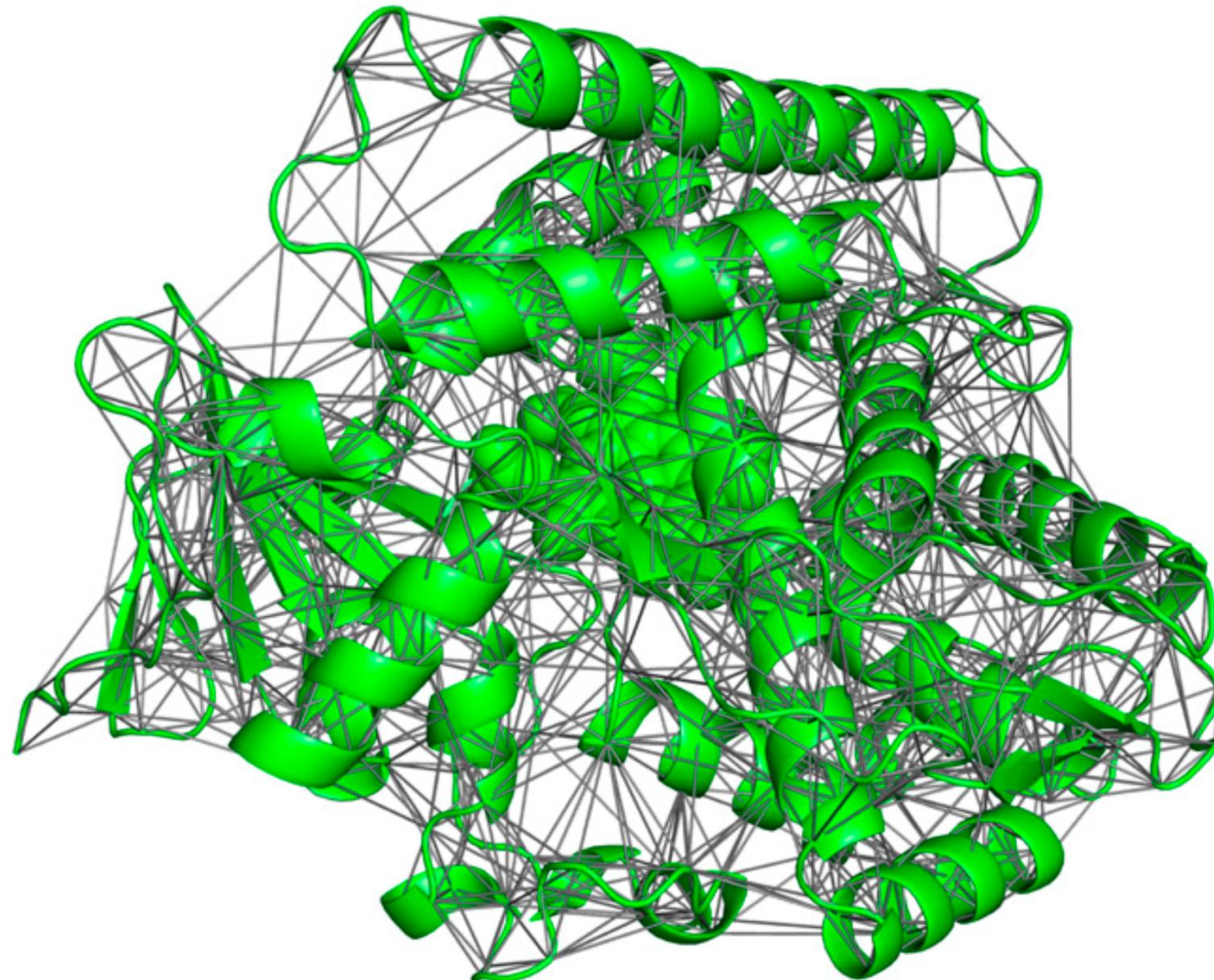


Example fits



Application: Designing P450s chimeras

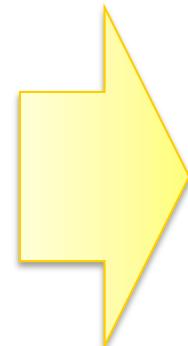
[with Phil Romero, Frances Arnold PNAS'13]



Design space

Parent
sequences

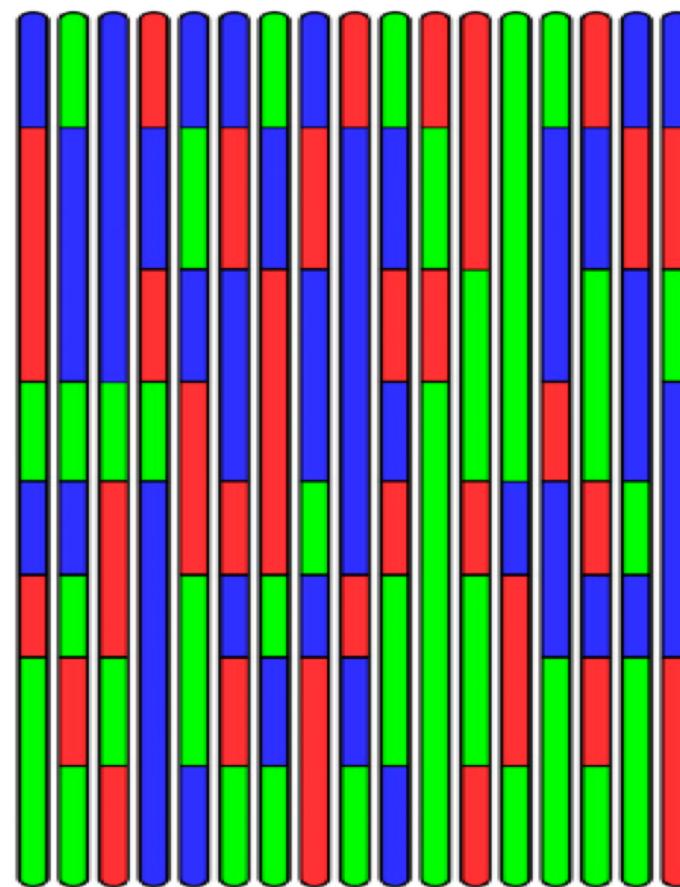
ABC



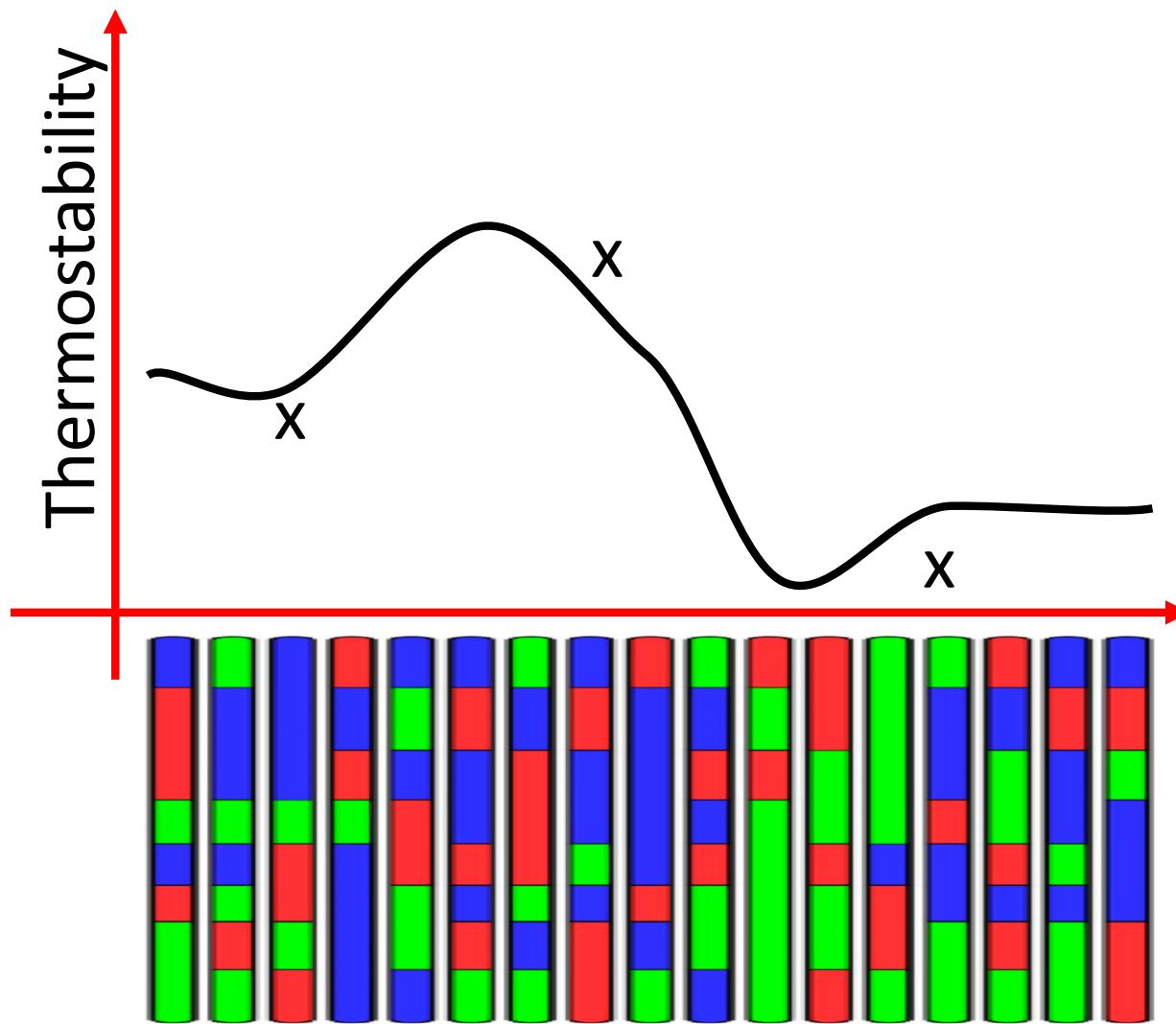
Candidate
designs

1 2 3 ...

n

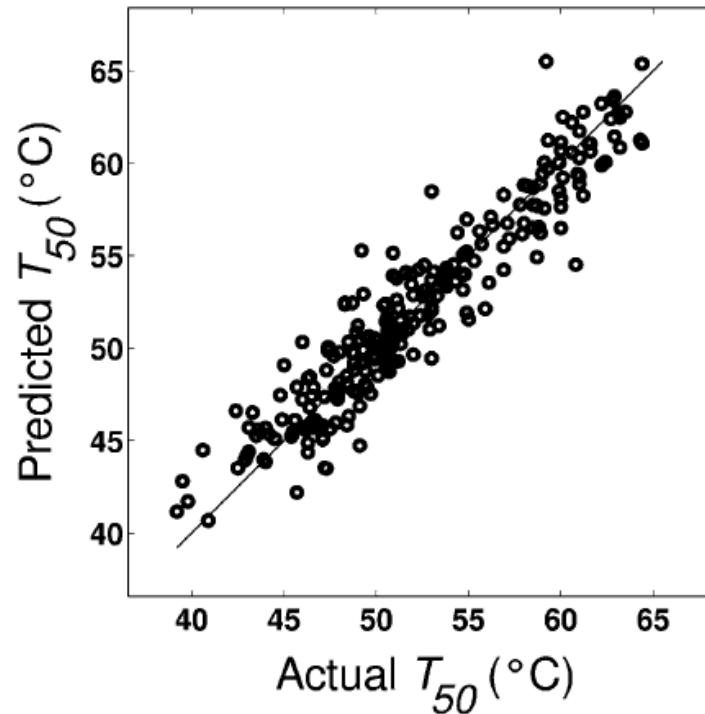
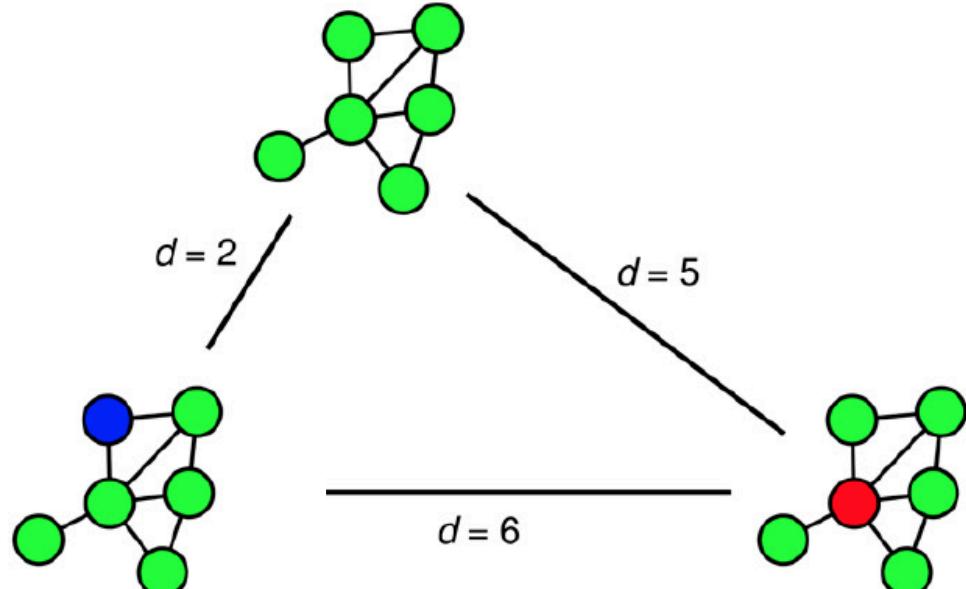


Protein Fitness Landscape



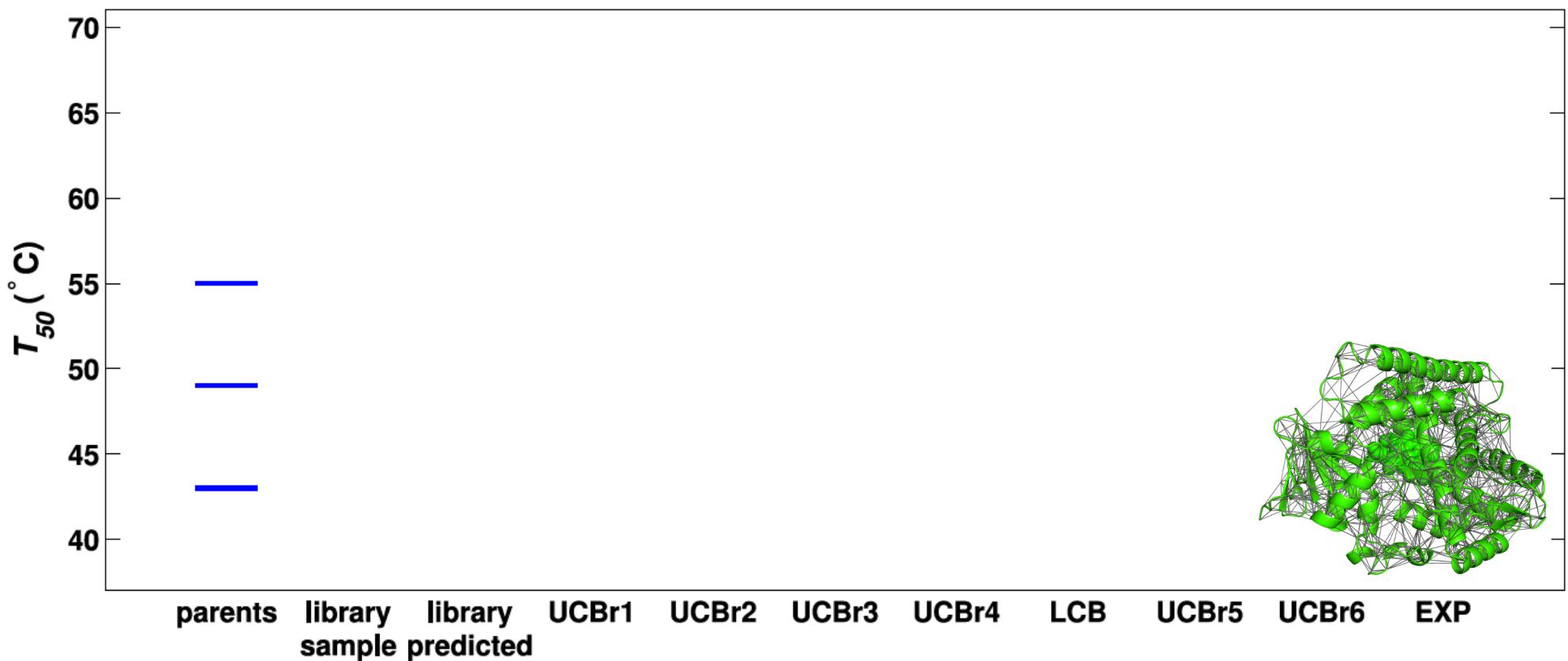
Application: Protein Engineering

[with Romero, Arnold, PNAS '13]



Wet-lab results

[w Romero, Arnold PNAS '13]



- Identification of new thermostable P450s chimera
- 5.3C more stable than best published sequence!**

Choosing kernels

- For a given kernel, how should we choose parameters?
 - Cross-validation! ☺
- How should we select suitable kernels?
 - Domain knowledge (dependent on data type)
 - «Brute force» (or heuristic) search
 - Use cross-validation
- **Learning kernels**
 - Much research on automatically selecting good kernels
(Multiple Kernel Learning; Hyperkernels; etc.)

Parameter demo

What about overfitting?

- Kernels map to (very) high-dimensional spaces.
- Why do we hope to be able to learn?
- **First attempt of an answer:**
(typically) # parameters \ll # dimensions. Why?
- Number of parameters = number of data points
(„non-parametric learning“)

What about overfitting?

- Kernels map to (very) high-dimensional spaces.
- Why do we hope to be able to learn?
- **Second attempt of an answer:**
- Overfitting can of course happen
(if we choose poor parameters)
- Can combat overfitting by regularization
 - This is already built into kernelized linear regression
(and SVMs), but **not** the kernelized Perceptron

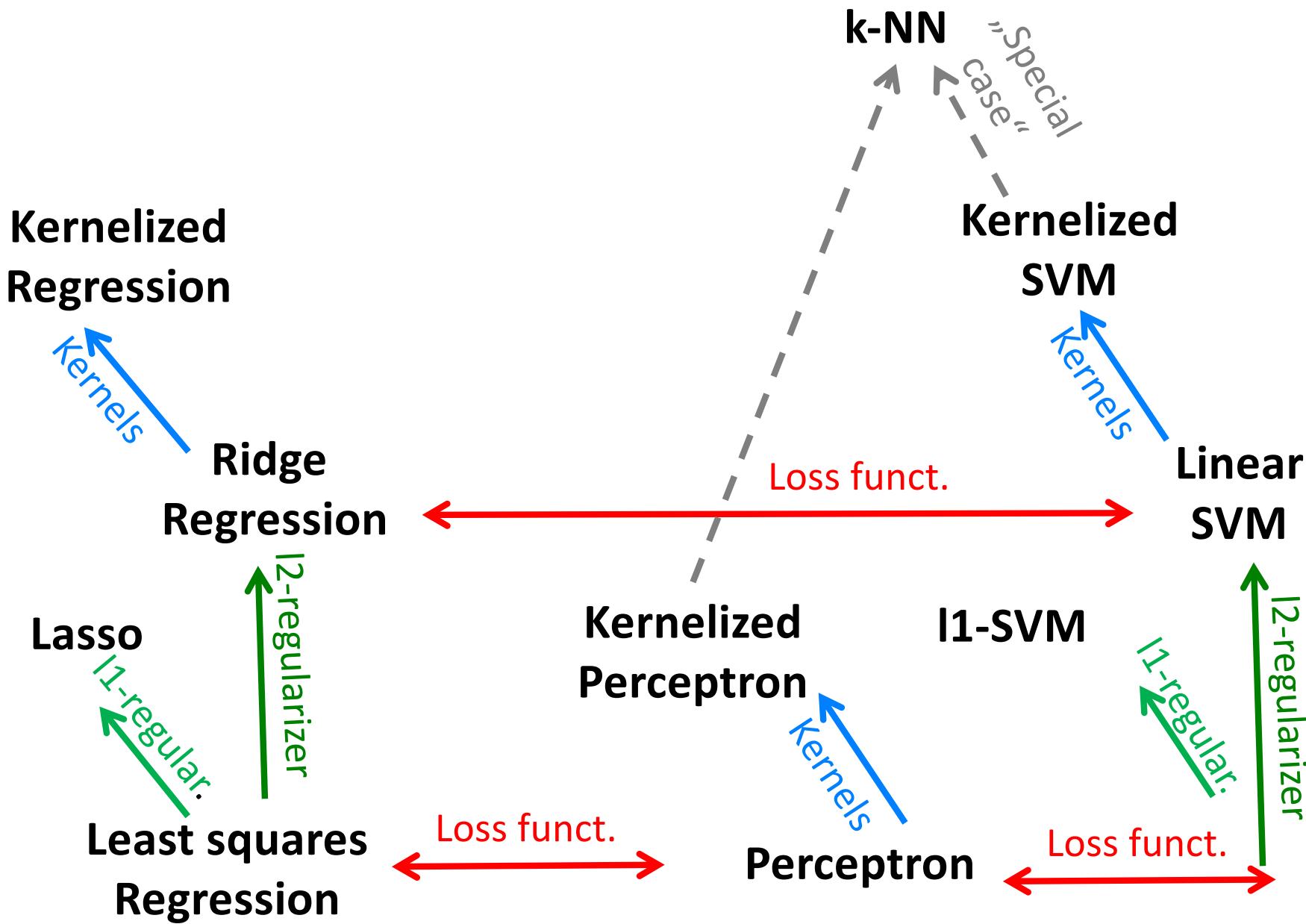
KLR: $\hat{\alpha} = \arg \min_{\alpha} \frac{1}{n} \|\alpha^T \mathbf{K} - \mathbf{y}\|_2^2 + \lambda \alpha^T \mathbf{K} \alpha$

SVM: $\hat{\alpha} = \arg \min_{\alpha} \frac{1}{n} \sum_{i=1}^n \max\{0, 1 - y_i \alpha^T \mathbf{k}_i\} + \lambda \alpha^T \mathbf{D_y K D_y} \alpha$

What you need to know

- Kernels are
 - (efficient, implicit) inner products
 - Positive (semi-)definite functions
 - Many examples (linear, polynomial, Gaussian/RBF, ...)
- The „Kernel trick“
 - Reformulate learning algorithm so that inner products appear
 - Replace inner products by kernels
- K-Nearest Neighbor classifier (and relation to Perceptron)
- How to choose kernels (kernel engineering etc.)
- **Applications:** Kernelized Perceptron / SVM; kernelized linear regression

Supervised learning big picture so far



Supervised learning summary so far

Representation/ features	Linear hypotheses; nonlinear hypotheses with nonlinear feature transforms; kernels		
Model/ objective:	Loss-function	+	Regularization
	Squared loss, 0/1 loss, Perceptron loss, Hinge loss		L^2 norm, L^1 norm
Method:	Exact solution, Gradient Descent, (mini-batch) SGD, Convex Programming, ...		
Evaluation metric:	Mean squared error, Accuracy		
Model selection:	K-fold Cross-Validation, Monte Carlo CV		