# Machine Learning from Data

Lecture 18: Spring 2021

# Today's Lecture
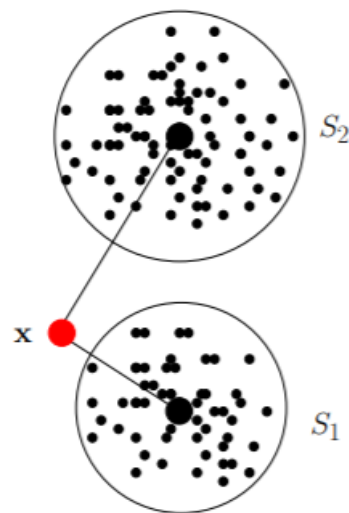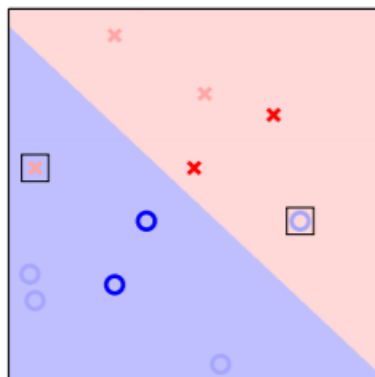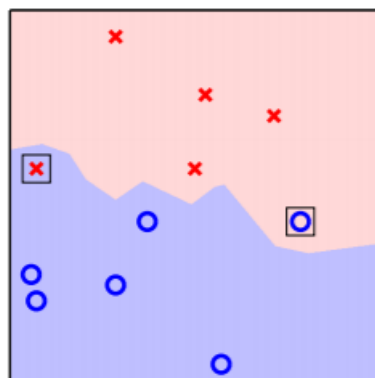
- Radial basis Functions  ( RBF )
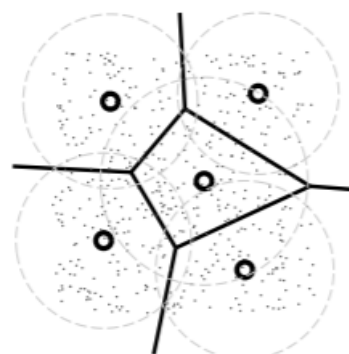- Non-Parametric RBF
- Parametric RBF
- K-RBF-Network

# Data Condensation and Nearest Neighbor Search

CNN

## Training Set Consistent



$S_2$

$S_1$

x

Branch and bound for finding nearest neighbors.

Lloyd's algorithm for finding a good clustering.

# Radial Basis Functions (RBF)

**$k$-Nearest Neighbor:** Only considers $k$-nearest neighbors.

    each neighbor has equal weight

What about using *all* data to compute $g(\mathbf{x})$?

**RBF:** Use all data.

    data further away from $\mathbf{x}$ have less weight.

# Non-parametric RBFs

## Regression

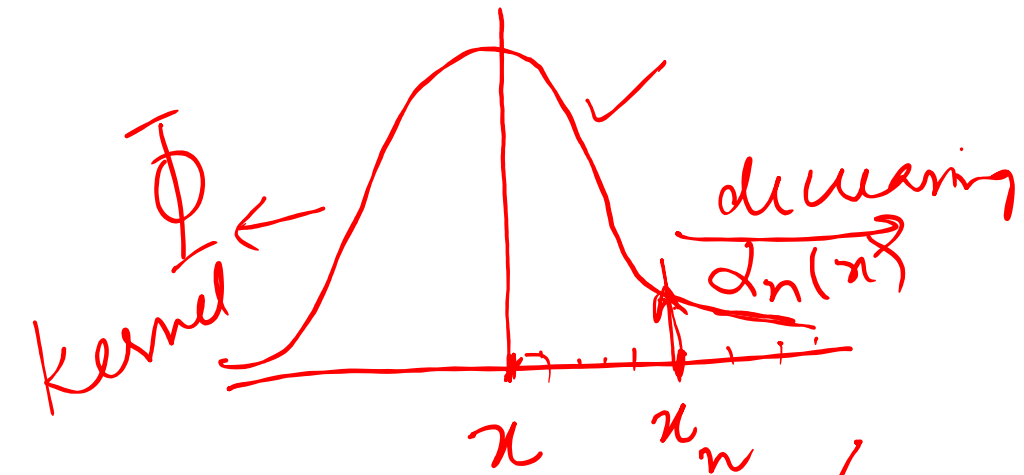$$KNN: \quad g(x) = \frac{1}{K} \sum_{i=1}^{K} y_{[i]}$$

$$or \quad g(x) = \sum_{i=1}^{K} \frac{1}{K} y_{[i]}$$

RBF

$$g(x) = \sum_{n=1}^{N} \alpha_n (x)$$

$$g(x) = \sum_{n=1}^{N} \frac{\alpha_n(x)}{\sum_{m=1}^{N} \alpha_m(x)} \cdot y_n$$

→ all the data

→ normalized to 1

$$\text{kernel} \leftarrow \Phi \quad \underline{\text{decreasing}} \atop \alpha_n(x)$$



$$x \qquad x_n$$

$$\therefore \ \alpha_n(x) = \Phi\left( \| x - x_n \| \right)$$

$$\gamma \rightarrow \text{scale} \atop \text{parameter}.$$

$$\alpha_n(x) = \Phi\left( \frac{\| x - x_n \|}{\gamma} \right) \underline{\text{Units}}$$

$$\Phi(s) = \frac{1}{1 + s^2} \qquad \left( \text{Decreasing} \atop \text{functions} \right)$$
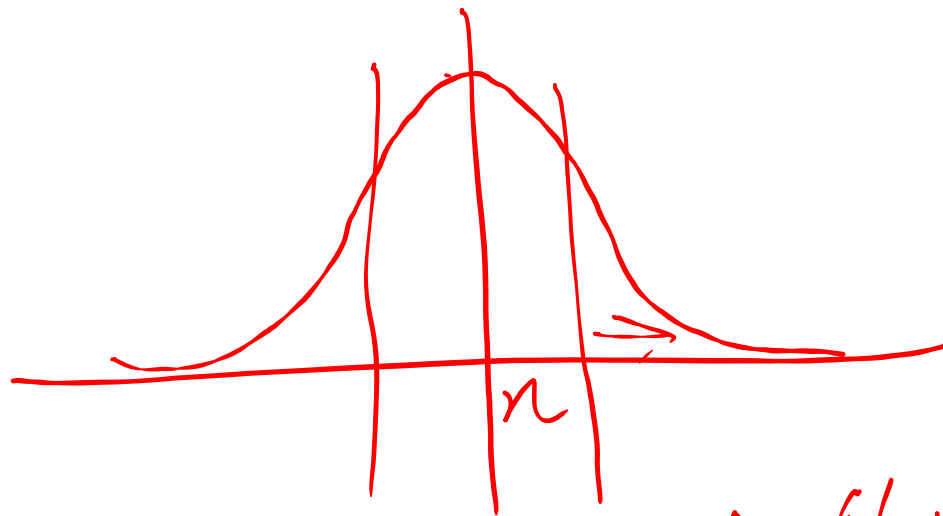
## 1) Window Kernel

$$\phi(s) = \begin{cases} 1 & \text{if } |s| \leq R \\ 0 & \text{if } |s| > R \end{cases}$$

## 2) Gaussian Kernel

$$\phi(s) = e^{-s^2/2}$$

$$g(n) = \sum_{n=1}^{N} \frac{\alpha_n(n)}{\sum_{m=1}^{M} \alpha_m(n)}$$

$$\alpha_n(x) = \phi\left(\left(\frac{\|x - x_n\|}{\gamma}\right)\right)$$

## Classification
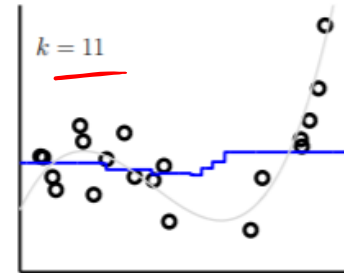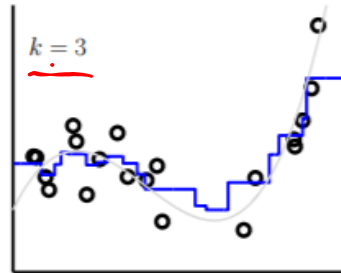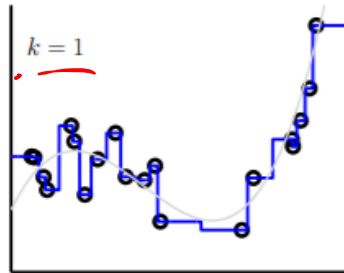
$$g(x) = \text{sign}\left( \sum_{n=1}^{N} \frac{\alpha_n(x)}{\sum_{m=1}^{M} \alpha_m(x)} \right) X$$

## Logistic Regression

$$g(x) = X \; [[y_n = +1]]$$
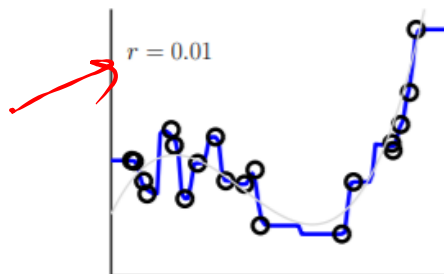
# Choice of Scale $r$

## Nearest Neighbor
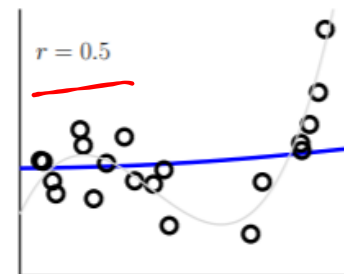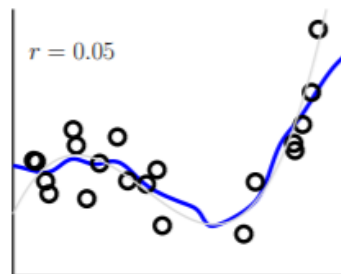


Choosing $k$:

$$k = 3$$

$$\boxed{k = \sqrt{N}}$$

CV

## Nonparametric RBF

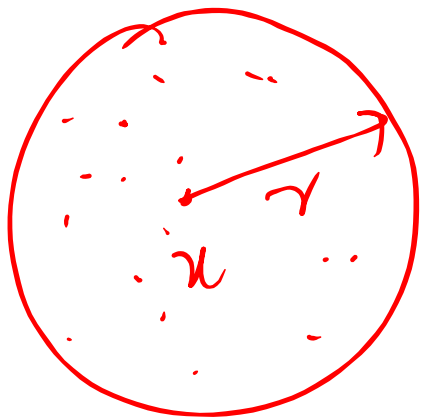

overfitting

underfitting

Choosing $r$:

$$r \sim \frac{1}{\sqrt[2d]{N}}$$

CV

$$\frac{||x - x_n||}{r} \to 0$$

$$k = \sqrt{N}$$



$$Vol \propto r^d$$

$$\# \text{ of data points} \propto N r^d$$

$$\cancel{N r^d} \qquad N r^d = \sqrt{N}$$

$$\boxed{r = \frac{1}{2\sqrt[d]{N}}} \quad (Recommendation)$$

# Highlights of Nonparametric RBF

1. Simple ('smooth' version of $k$-NN rule).

2. No training.

3. Near optimal $E_{\text{out}}$. $\longleftarrow$ $\frac{k}{N} \to 0$

4. Easy to justify classification to customer.

} A **good!** method

5. Can do classification, multi-class, regression, logistic regression.

6. **Computationally demanding**.

$r = 0.01$

$e^{-s/2}$

$v = 0.3$

$$g(x) = \sum_{n=1}^{N} \left[ \frac{\alpha_n(x)}{\sum_{m=1}^{M} \alpha_m(x)} \right] y_n$$

# PARAMETRIC RBF

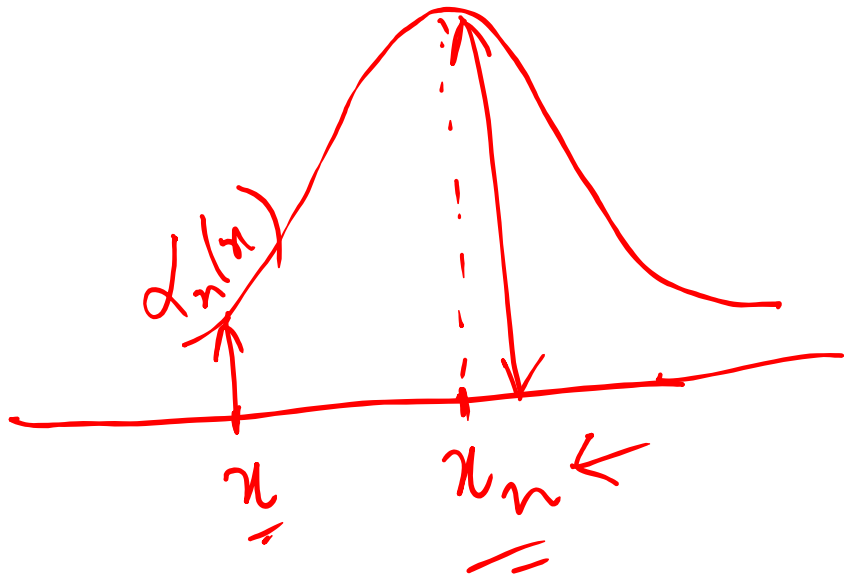$$g(x) = \sum_{n=1}^{N} \underbrace{\frac{y_n}{\sum_{m=1}^{M} \alpha_m(x)}}_{\omega_n(x)} \cdot \phi\left(\frac{\|x - x_n\|}{r}\right) \rightarrow \text{symmetric}$$

Earlier $x$ as the center



$$g(x) = \sum_{n=1}^{N} \omega_n(x) \, \phi\left(\frac{\|x - x_n\|}{r}\right)$$

Chasing curve $\longrightarrow$ serving data point.

$w_m^{(n)}$ $v$ $w_n^{(n)}$

$x_m$ $x$ $x_n$

$g(x)$

Approximation

$$g(x) = \sum_{n=1}^{N} w_n \, \phi\left( \frac{\|x - x_n\|}{r} \right)$$

PARAMETRIC RBF.

$r = 0.1$

$r = 0.3$

$$g(x) = \sum_{n=1}^{N} \omega_n \, \phi\left(\frac{\|x - x_n\|}{\gamma}\right)$$

parameters

we want

$$g(x_1) = \sum_{n=1}^{N} \omega_n \, \phi\left(\frac{\|x_1 - x_n\|}{\gamma}\right) = y_1$$

$$\vdots$$

$$g(x_N) = \sum_{n=1}^{N} \omega_n \, \phi\left(\frac{\|x_N - x_n\|}{\gamma}\right) = y_N$$

N eqns x N unknowns.

$$x \longrightarrow Z \quad (\underline{N} \text{ dimensional} \quad \text{feature transform}).$$

$$\begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_N \end{bmatrix} = \begin{bmatrix} \phi\left(\dfrac{\|x - x_1\|}{\gamma}\right) \\ \vdots \\ \phi\left(\dfrac{\|x - x_N\|}{\gamma}\right) \end{bmatrix}$$

Similarity features!

$\longrightarrow$ Depend on the data.

$$Z = \begin{bmatrix} - z_1^T - \\ - z_2^T - \\ \vdots \\ - z_N^T - \end{bmatrix} \qquad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

data matrix $\longleftarrow$ $Zw \approx y \longrightarrow$ target vector.

## Regression

$$w_{lin} = Z^\dagger y = \left( Z^T Z \right)^{-1} Z^T y \qquad \longrightarrow N \times N \text{ matrix}$$

Classification $\longrightarrow$ PLA $\longrightarrow$ Pocket Algo.

Logistic Regression $\longrightarrow$ Gradient Descent.
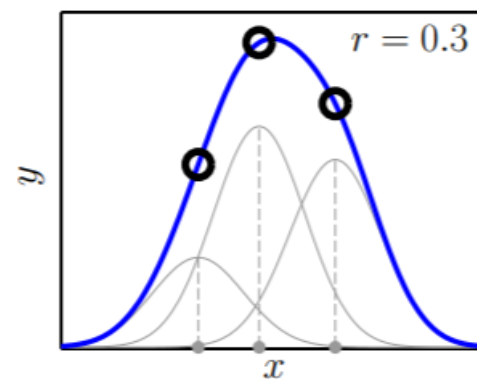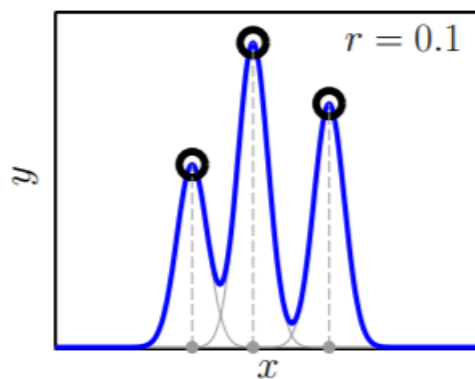
Lin. Model with similarity feature.

$$h(\mathbf{x}) = \sum_{n=1}^{N} w_n \cdot \phi\left(\frac{\|\mathbf{x} - \mathbf{x}_n\|}{r}\right) = \mathbf{w}^\mathsf{T}\mathbf{z}$$

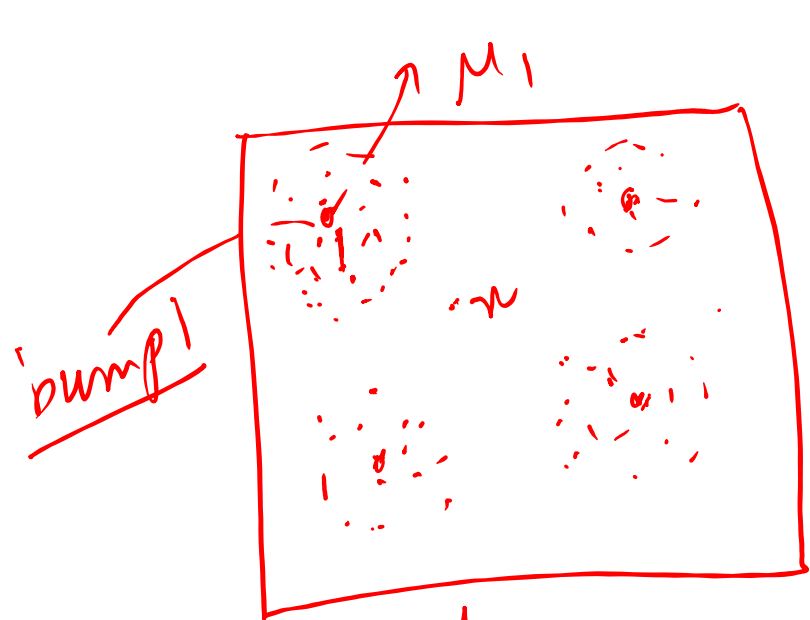$$\mathbf{z} = \boldsymbol{\Phi}(\mathbf{x}) = \begin{bmatrix} \phi_1(\mathbf{x}) \\ \phi_2(\mathbf{x}) \\ \vdots \\ \phi_N(\mathbf{x}) \end{bmatrix}, \quad \phi_n(\mathbf{x}) = \phi\left(\frac{\|\mathbf{x} - \mathbf{x}_n\|}{r}\right). \qquad \mathbf{Z} = \begin{bmatrix} -\ \mathbf{z}_1^\mathsf{T}\ - \\ -\ \mathbf{z}_2^\mathsf{T}\ - \\ \vdots \\ -\ \mathbf{z}_N^\mathsf{T}\ - \end{bmatrix} = \begin{bmatrix} -\ \boldsymbol{\Phi}(\mathbf{x}_1)^\mathsf{T}\ - \\ -\ \boldsymbol{\Phi}(\mathbf{x}_2)^\mathsf{T}\ - \\ \vdots \\ -\ \boldsymbol{\Phi}(\mathbf{x}_N)^\mathsf{T}\ - \end{bmatrix}$$

Fit the data $(h(\mathbf{x}_n) = y_n)$:

$$\mathbf{w} = \mathbf{Z}^\dagger \mathbf{y} = (\mathbf{Z}^\mathsf{T}\mathbf{Z})^{-1}\mathbf{Z}^\mathsf{T}\mathbf{y}$$

$\mu_1$

'bump'

2-d

4 bumps/peaks

N peaks

Reduce the no. of parameter.

$\rightarrow$ Decide # of bumps ... K with centus $\{\mu_1, \mu_2 ... \mu_k\}$

$\rightarrow$ Center must cover the data.

$$g(x) = \omega_0 + \sum_{j=1}^{K} \omega_j \phi\left(\frac{||x-\mu_j||}{r}\right)$$

$$x \rightarrow Z = \begin{bmatrix} \phi\left(\frac{||x-\mu_1||}{r}\right) \\ \phi\left(\frac{||x-\mu_k||}{r}\right) \end{bmatrix}$$

$\checkmark$ Reduced

$\omega^T Z$

$\rightarrow$ After selecting $\mu_j$'s

$$g(x) = \underline{\underline{\omega_0}} + \sum_{j=1}^{K} \omega_j \cdot \phi\left(\underbrace{\frac{\|x - \mu_j\|}{\gamma}}_{\text{similarity}}\right)$$



Radial Basis function Network

$\mu$'s ?

i) Find <u>centers</u> of clusters in the data.

    Lloy'd's Algorithm. ✓

    Not using y values.

    → Unsupervised step.

ii) You have a linear model with $\omega$'s

    PLA, Pseudo Inverse, grad descent

Picking $r$ values. ⟶ $k$ bumps in $d$ dimensions.

Vol. of data covered $\sim r^d$     $r \propto \dfrac{1}{\sqrt[d]{K}}$

$K$ bumps $\sim K r^d$ ⟵

<u>Recommendation</u>    $K r^d \sim R$

                   $r \sim \dfrac{R}{\sqrt[d]{K}}$

# Fitting the Data

**Fitting the RBF-network to the data (given $k, r$):**

1: Use the inputs X to determine $k$ centers $\mu_1, \ldots, \mu_k$.

2: Compute the $N \times (k+1)$ feature matrix Z

$$Z = \begin{bmatrix} - & \mathbf{z}_1^T & - \\ - & \mathbf{z}_2^T & - \\ & \vdots & \\ - & \mathbf{z}_N^T & - \end{bmatrix} = \begin{bmatrix} - & \boldsymbol{\Phi}(\mathbf{x}_1)^T & - \\ - & \boldsymbol{\Phi}(\mathbf{x}_2)^T & - \\ & \vdots & \\ - & \boldsymbol{\Phi}(\mathbf{x}_N)^T & - \end{bmatrix}, \text{ where } \boldsymbol{\Phi}(\mathbf{x}) = \begin{bmatrix} 1 \\ \phi_1(\mathbf{x}) \\ \vdots \\ \phi_k(\mathbf{x}) \end{bmatrix}, \phi_j(\mathbf{x}) = \phi\left(\frac{|\mathbf{x}-\mu_j|}{r}\right)$$

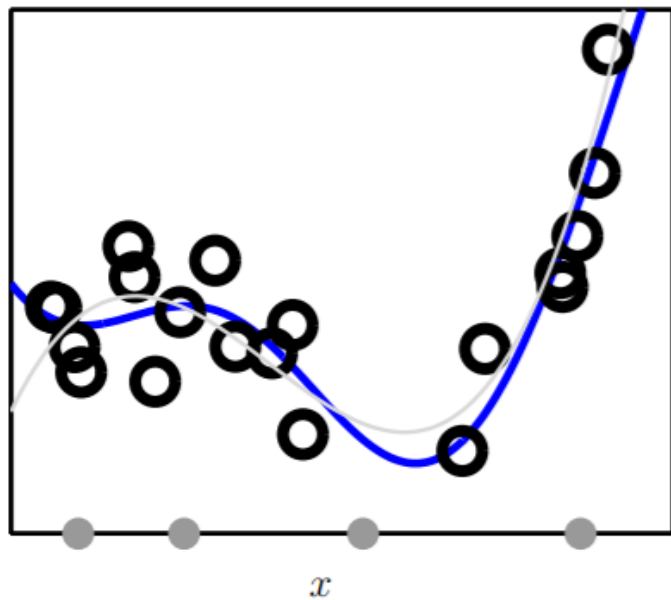Each row of Z is the RBF-feature corresponding to $\mathbf{x}_n$ (with dummy bias coordinate 1).

3: Fit the *linear* model $Z\mathbf{w}$ to $\mathbf{y}$ to determine the weights $\mathbf{w}^*$.
   classification: PLA, pocket, linear programming,....
   regression: pseudoinverse.
   logistic regression: gradient descent on cross entropy error.
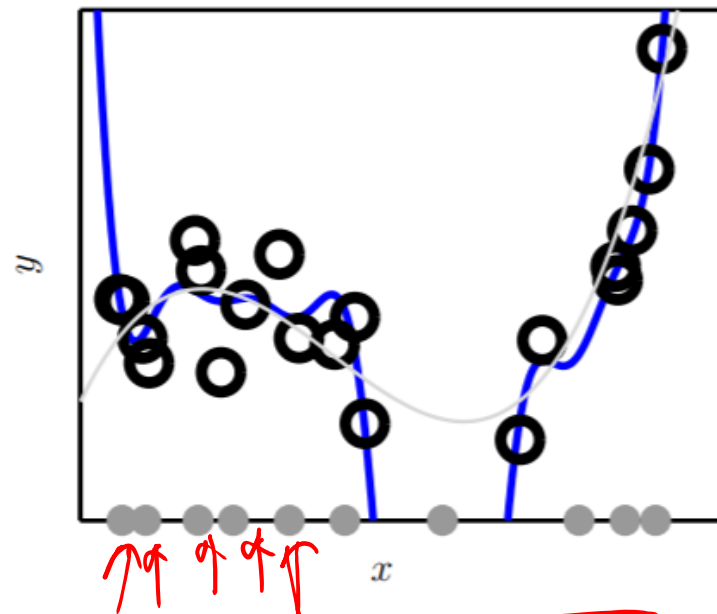
$k = 4, \ r = \frac{1}{k}$

$k = 10, \ r = \frac{1}{k}$

CROSS VALIDATION

Reasonable

white $\longrightarrow$ target fn.

$\mathbf{w} = (\mathbf{Z}^{\mathrm{T}}\mathbf{Z})^{-1}\mathbf{Z}^{\mathrm{T}}\mathbf{y}$
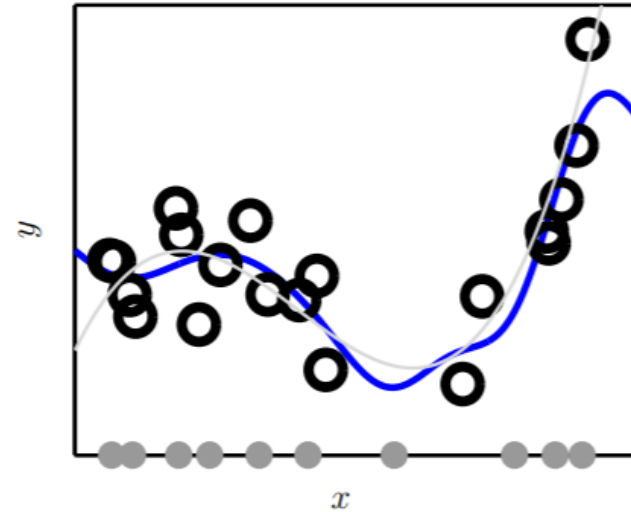
Linear Model

# Use Regularization to Fight Overfitting

$$k = 10, \ r = \frac{1}{k}$$
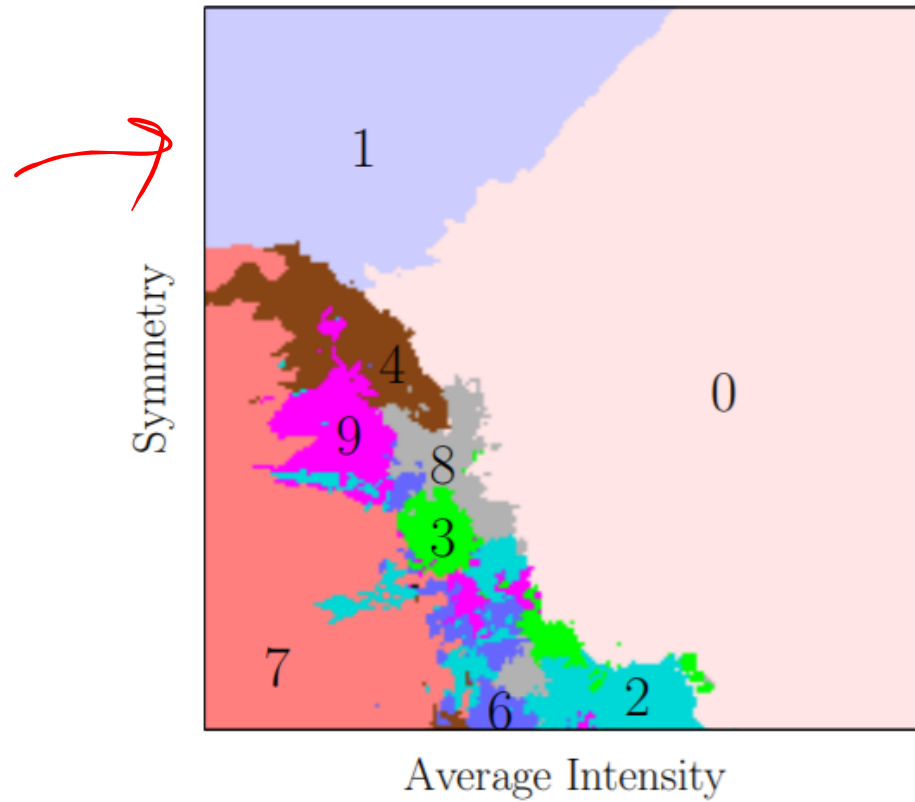
$$k = 10, \ r = \frac{1}{k}, \text{ regularized}$$

*Overfitting*

$$\mathbf{w} = (\mathbf{Z}^{\mathrm{T}}\mathbf{Z} + \lambda\mathbf{I})^{-1}\mathbf{Z}^{\mathrm{T}}\mathbf{y}$$
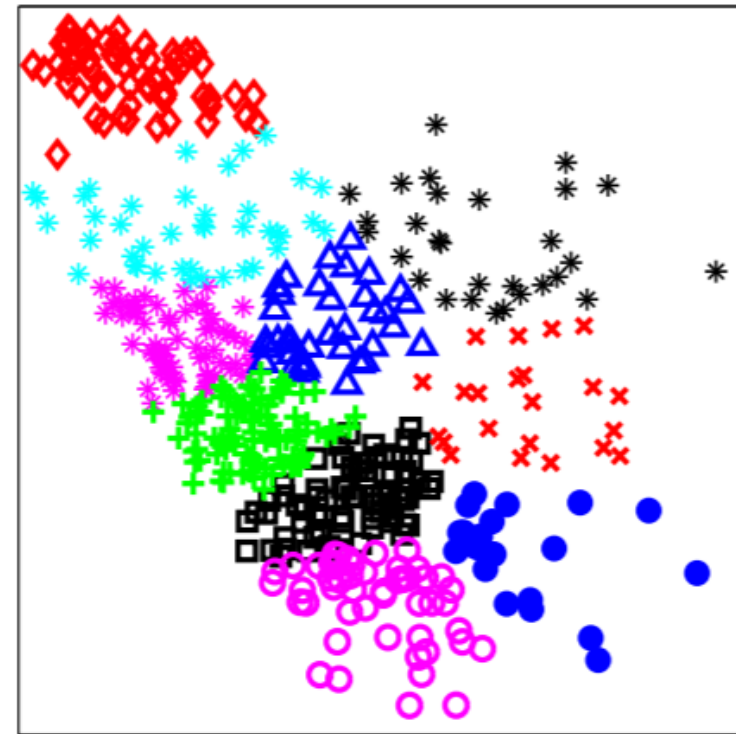
# Reflecting on the $k$-RBF-Network

1. We derived it as a 'soft' generalization of $k$-NN rule.

      Can also be derived from regularization theory.

      Can also be derived from noisy interpolation theory.

2. Can use nonparametric or parametric versions.

3. Given centers, 'easy' to learn the weights using techniques from linear models.

      A linear model with an adaptable nonlinear transform.

4. We used uniform bumps – can have different shapes $\Sigma_j$.

5. **NEXT:** How to better choose the centers: unsupervised learning.

# A Peek at Unsupervised Learning

10 Clustering of Data

# Thanks!