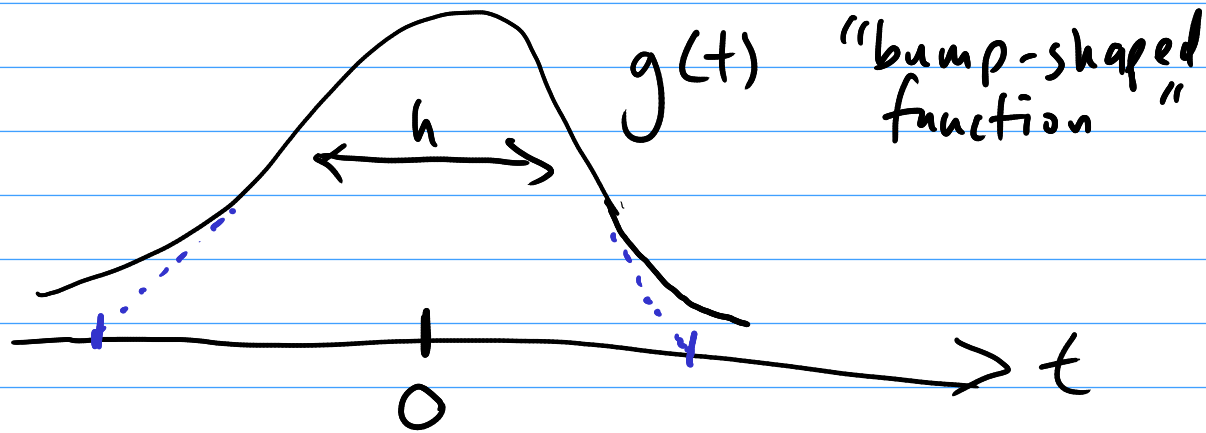Detailed project instructions coming soon!

# Radial basis function networks

Goals: introduce RBFs
fitting w/ fixed centers
differences from kernel smoothers

# Kernel smoother

$h = \text{bandwidth}$



$g(t)$ "bump-shaped function"

$$g(t) = \exp\left(\frac{-t^2}{2h^2}\right)$$

Gaussian function

$$\underbrace{K(\vec{x}, \vec{x}')}_{\text{kernel}} = g\left(\underbrace{\|\vec{x} - \vec{x}'\|}_{\text{RBF}}\right)$$

$$f(\vec{x}) = \sum_{i=1}^{n} w_i \, y_i \, K(\vec{x}, \vec{x}_i)$$

Nadaraya-Watson

$$= \frac{\sum_{i=1}^{n} y_i \, K(\vec{x}, \vec{x}_i)}{\sum_{i=1}^{n} K(\vec{x}, \vec{x}_i)} \quad \right\} \text{no fitting except } h$$

- gives smooth $f$ because $g$ smooth

- biased depending on density of training samples

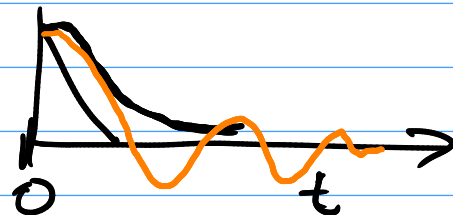# Radial basis function (RBF) Kernel function

w/ $k(\vec{x}, \vec{x}') = g(\|\vec{x} - \vec{x}'\|)$

kernel function

$\underbrace{\quad}$ similarity of $\vec{x}, \vec{x}'$

$\underbrace{\quad}$ depends only on distance

RBF

RBF network has the form

$$f(\vec{x}) = \sum_{i=1}^{m} w_i \, g(\|\vec{x} - \vec{c}_i\|)$$

$m$ — number of centers

$w_i$ — weights

$\vec{c}_i$ — centers

Differences:
- no $y_i$, learn from data
- centers are not on data pts (in general)
- # centers could be different than # data pts
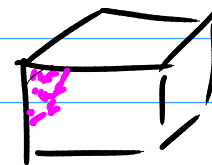
# Simple way to pick centers

1) Set $m < n$     #training pts $\downarrow$

Sample $I_j \sim \text{Uniform}(\{1, \dots, n\})$

Set $\vec{c}_j = \vec{x}_{I_j}$

Roughly centers distributed like the data

2) Form a grid over $\vec{x}$ domain of $m$ points



Having $m < n$ could increase efficiency a lot

# Fitting an RBF network

$\mathcal{X}$ = training set $\qquad (\vec{x}_i, y_i)_{i=1}^{n} \qquad$ regression or classification

**Fixed centers:** $\quad \min_{\vec{w}} \; \sum_{i=1}^{n} \ell\left(f(\vec{x}_i), y_i\right) + R(\vec{w})$

$$f(\vec{x}_i) = \sum_{j=1}^{m} w_j \underbrace{g\left(\|\vec{x}_i - \vec{c}_j\|\right)}_{G_{ij}}$$

$$= \sum_{j=1}^{m} G_{ij} w_j$$

$$= \left(G \vec{w}\right)_i$$

$$G = \overbrace{\begin{bmatrix} g\left(\|\vec{x}_1 - \vec{c}_1\|\right) & g\left(\|\vec{x}_1 - \vec{c}_2\|\right) \cdots \\ g\left(\|\vec{x}_2 - \vec{c}_1\|\right) & \ddots \\ \vdots & \ddots \\ g\left(\|\vec{x}_n - \vec{c}_1\|\right) & \end{bmatrix}}^{m \text{ centers}}$$
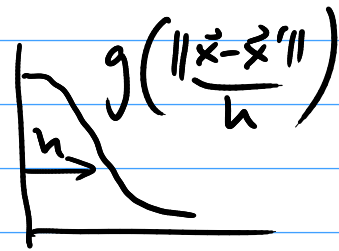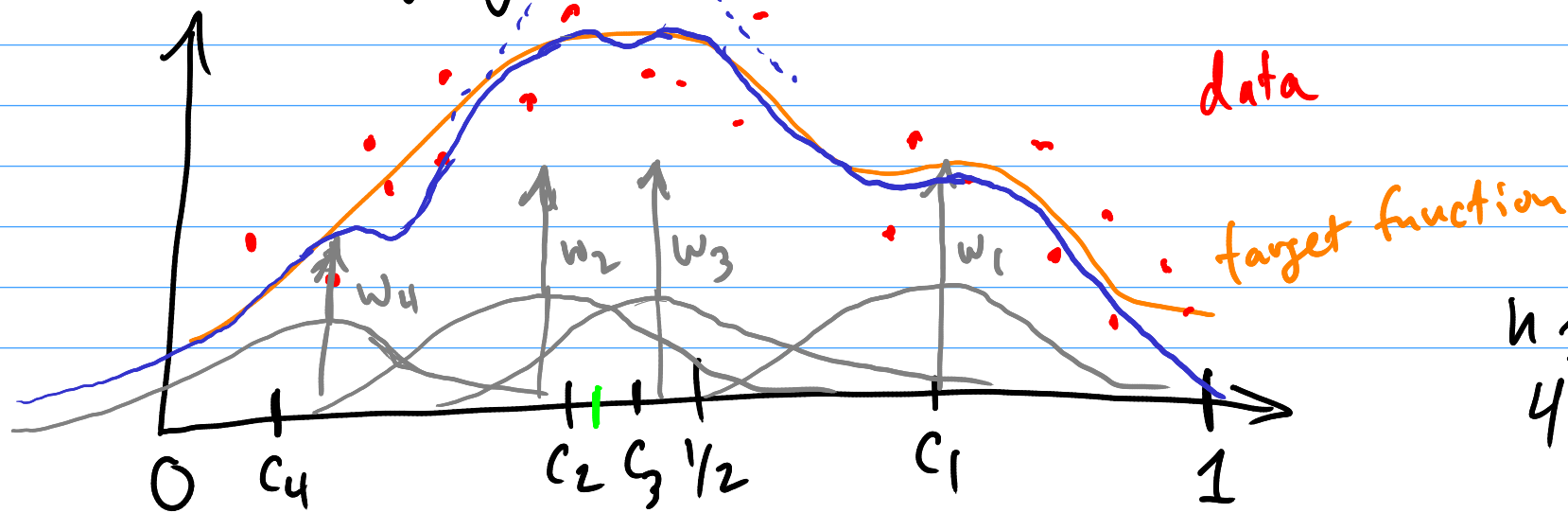
matrix of similarities
$n \times m$

$$f(\vec{z}_i) = (G\vec{w})_i \implies \hat{y} = G\vec{w}$$

$\implies$ fitting $\vec{w}$ just like fitting a linear model

$\implies$ If $\ell$, $R$ convex, then finding $\vec{w}$ is a convex optimization

If you want to fit $\vec{c}_i$, nonconvex, can still try gradient descent.



data

target function

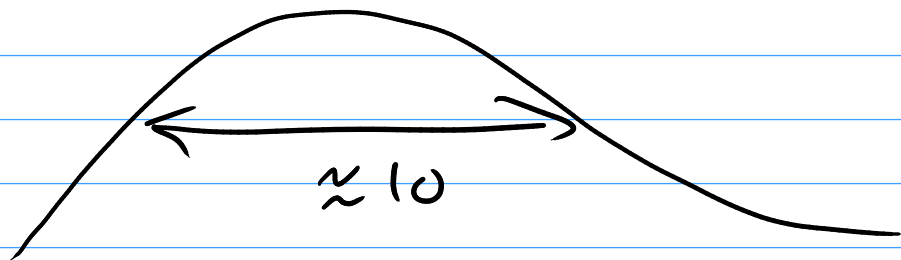$g\left(\frac{\|\vec{x} - \vec{x}'\|}{h}\right)$

$h \approx 0.5$
4 centers

# Practical considerations:

- bandwidth $h$ very important. Depends on
  - smoothness of target
  - # of data pts $n$, # of centers
  - input dim $d$
- $m$ also important
  - might pick $m$ based on compute

*use C.V. to pick*



$\approx 1$    $\approx 10$

$h=1$         $h=10$

$$\min_{\beta_0, \vec{\beta}} \quad \| X\vec{\beta} + \mathbb{1}\beta_0 - y \|^2 + \lambda \| \vec{\beta} \|_1$$

Cost

$$C(\vec{\beta}, \beta_0) = \sum_{i=1}^{n} \left( \vec{x}_i^T \vec{\beta} + \beta_0 - y_i \right)^2 + \lambda \sum_{i=1}^{d} |\beta_i|$$

$$\frac{\partial C}{\partial w_i} = 0 \qquad C = \| X\vec{w} - \vec{y} \|^2, \quad \nabla_{\vec{w}} C = 2X^T X \vec{w} - 2X^T \vec{y}$$

$i^{\underline{th}}$ element
of gradient $(\nabla_{\vec{w}} C)_i$

Pick out entry $i$
write as summation