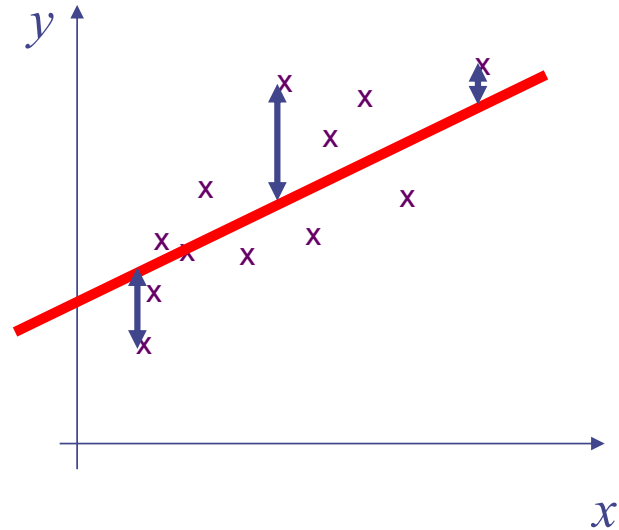


Machine Learning

4771

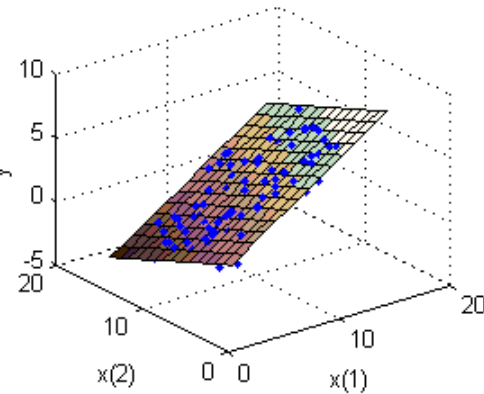
Instructor: Itsik Pe'er

Reminder: Regression

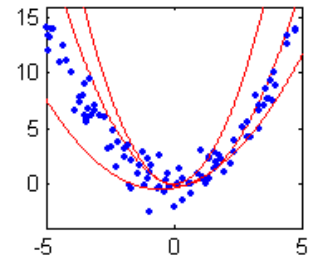


- 1D Linear:
 - Loss
 - Empirical risk
 - Least-squares

- Multi-D Linear: matrix form

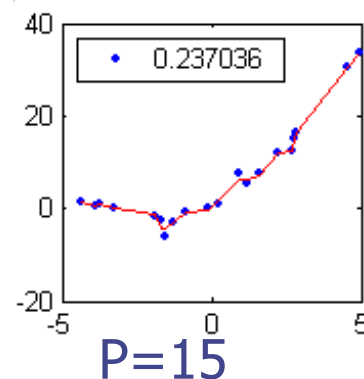
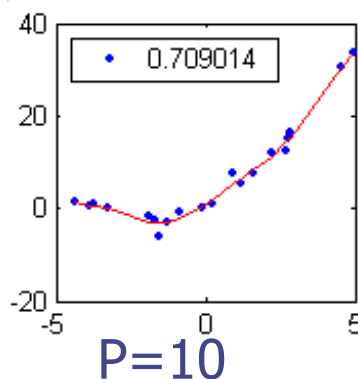
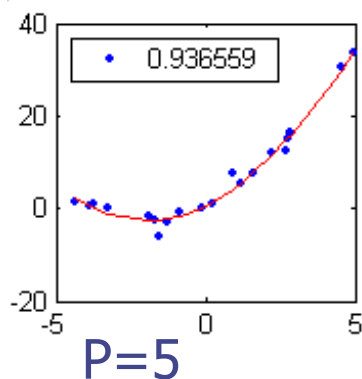
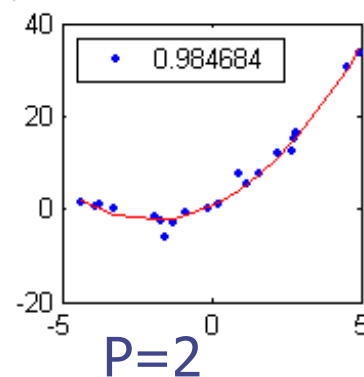
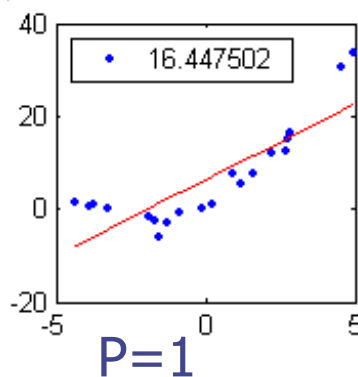
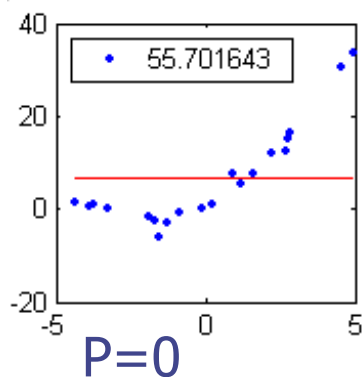


- Polynomial



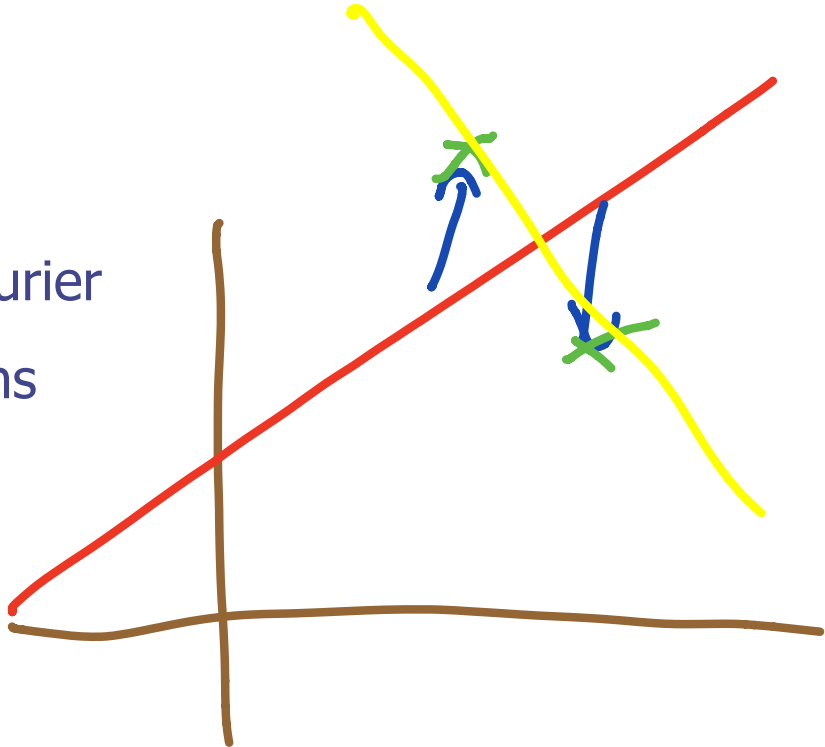
Underfitting/Overfitting

- Try varying P . Higher P fits a more complex function class
- Observe $R(\theta^*)$ drops with bigger P



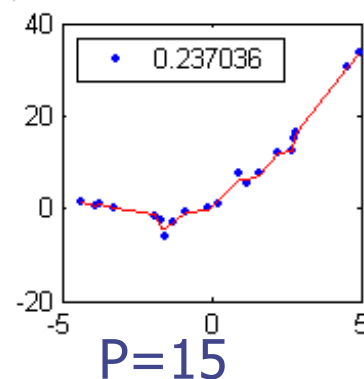
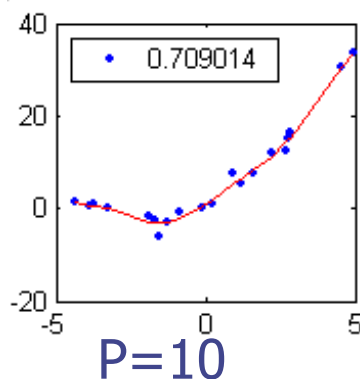
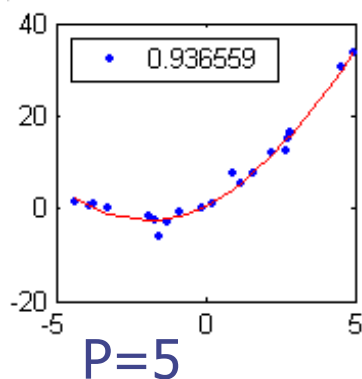
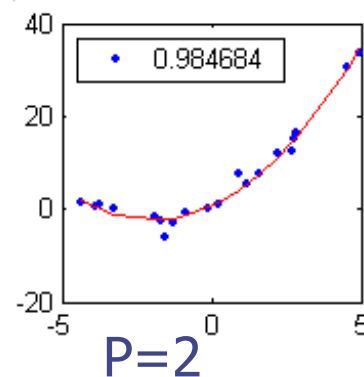
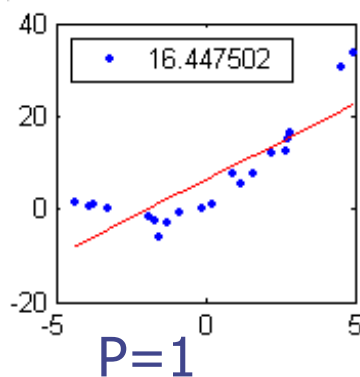
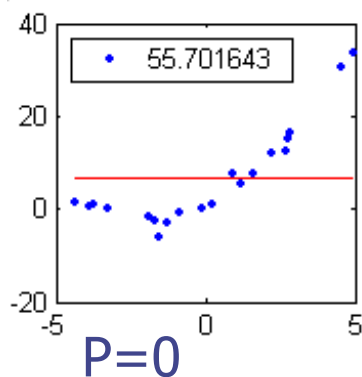
Class 4

- Overfitting
- Additive models: Fourier
- Radial basis functions



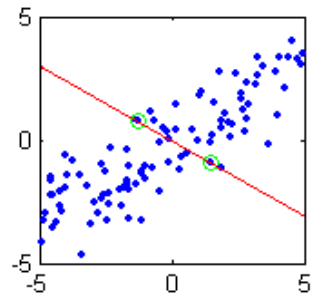
Underfitting/Overfitting

- Try varying P . Higher P fits a more complex function class
- Observe $R(\theta^*)$ drops with bigger P



Evaluating The Regression

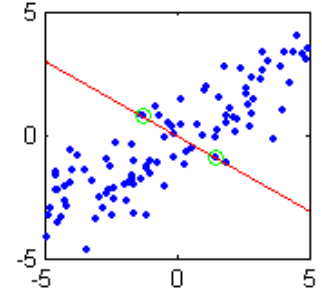
- Unfair to use empirical to find best order P
- High P (vs. N) can overfit, even linear case!
- $\min R(\theta^*)$ not on training but on future data
- Want model to *Generalize* to future data



$$R(\theta) = \int (y - f(x, \theta))^2 p(x, y) dx dy$$

Evaluating The Regression

- Unfair to use empirical to find best order P
- High P (vs. N) can overfit, even linear case!
- $\min R(\theta^*)$ not on training but on future data
- Want model to *Generalize* to future data



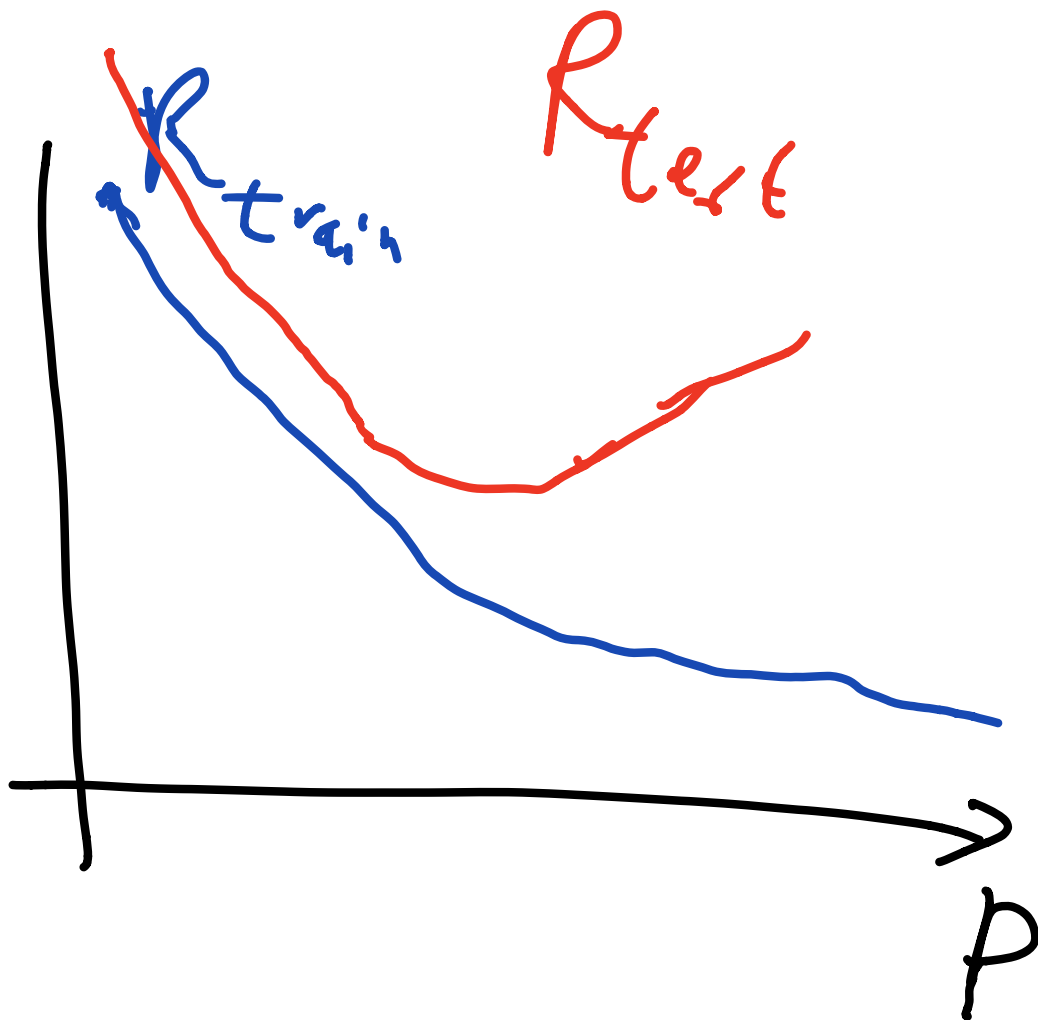
True loss: $R_{true}(\theta) = \int p(x, y) \frac{1}{2} (y - \theta^T x)^2 dx dy$

- One approach: split data into training / testing portion

$$\{(x_1, y_1), \dots, (x_N, y_N)\} \quad \{(x_{N+1}, y_{N+1}), \dots, (x_{N+M}, y_{N+M})\}$$

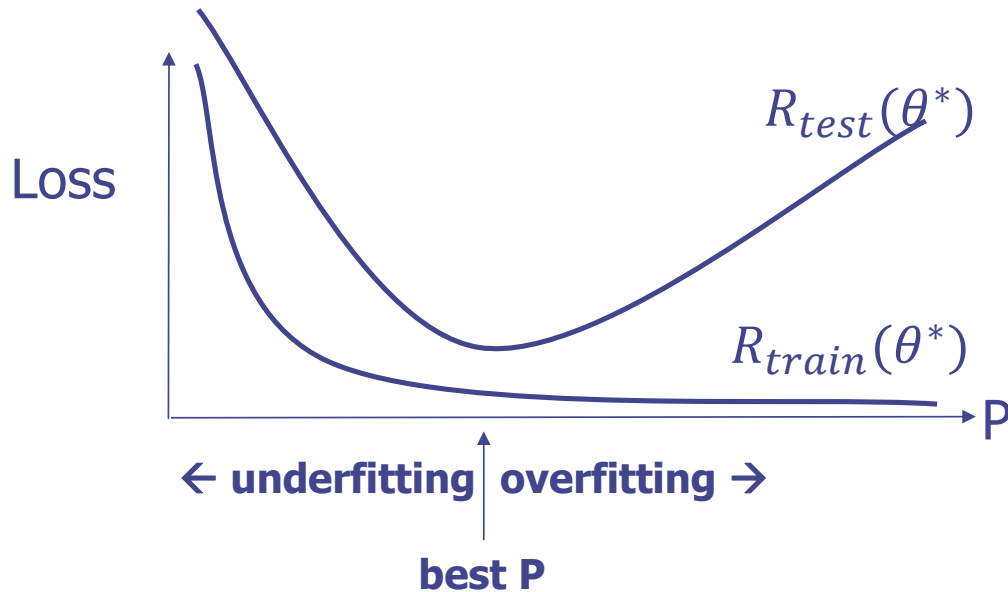
- Estimate θ^* with training loss: $R_{train}(\theta) = \frac{1}{2N} \sum_{i=1}^N (y_i - \theta^T x_i)^2$

- Evaluate P with testing loss: $R_{test}(\theta^*) = \frac{1}{2M} \sum_{i=N+1}^{N+M} (y_i - \theta^{*T} x_i)^2$



Crossvalidation

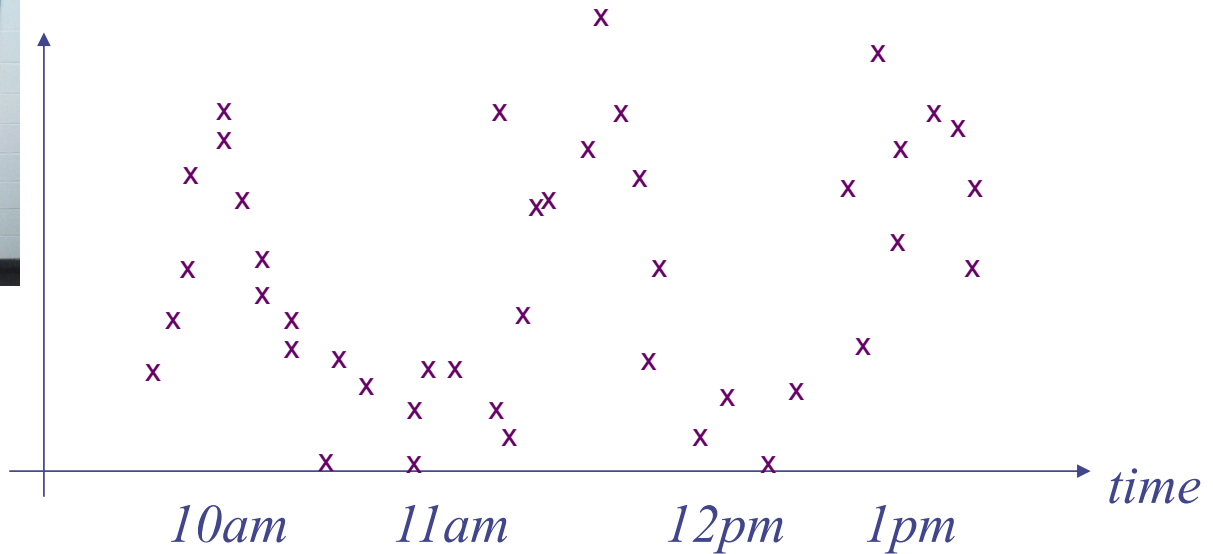
- Try fitting with different polynomial order P
- Select P which gives lowest $R_{test}(\theta^*)$



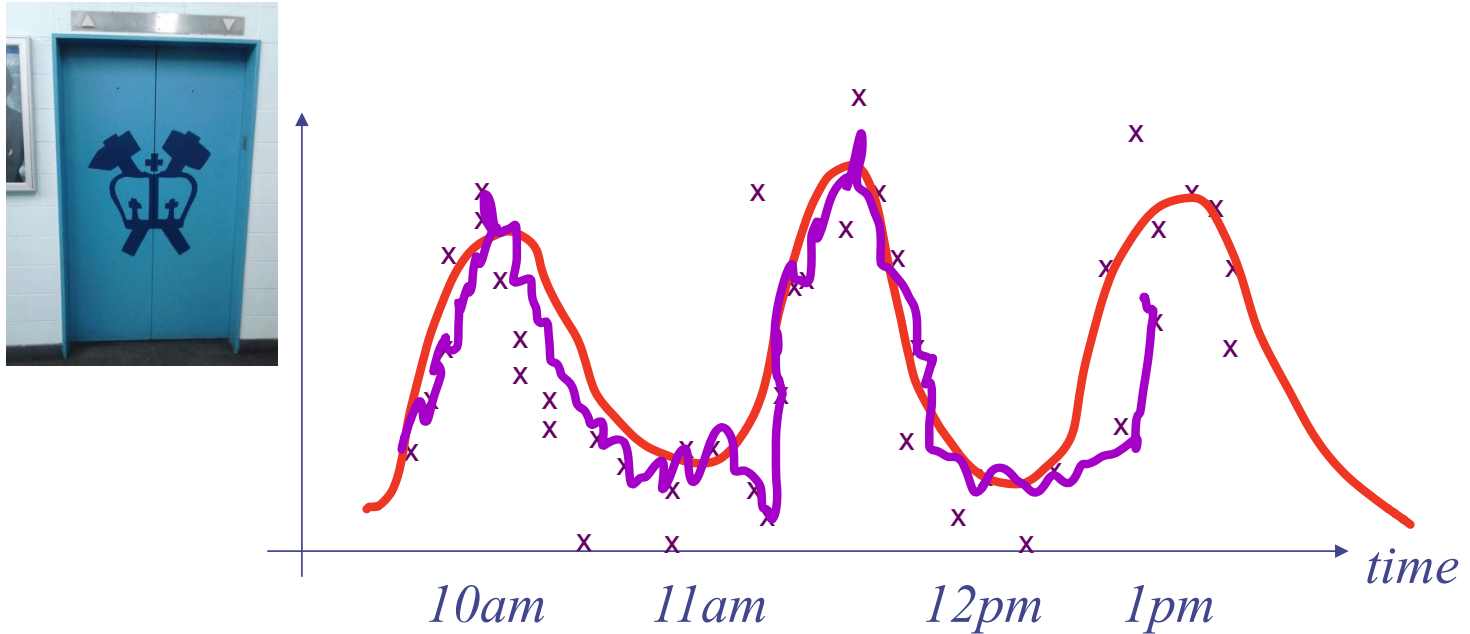
- Think of P as a measure of the complexity of the model
- Higher order polynomials are more flexible and complex



Example: Temporal data



Example: Temporal data



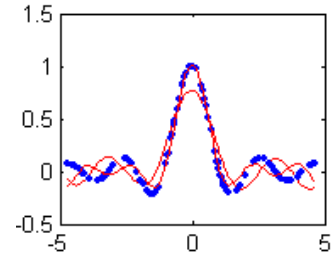
- ◆ Need to fit periodic behavior
- ◆ Cycle: 90min, daily, weekly, annual

Sinusoidal Basis Functions

- General functions, not just polynomials:

$$f(x; \theta) = \sum_{p=1} \theta_p \phi_p(x) + \theta_0$$

- These are generally called **Additive Models**
- Regression adds linear combinations of the basis fn's



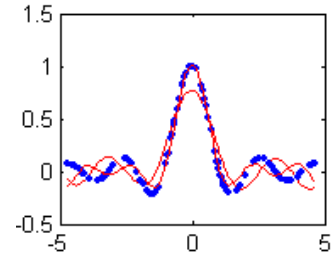
$$\phi_{2p}(x) = \sin px$$

$$\phi_{2p+1}(x) = \cos px$$

Sinusoidal Basis Functions

- General functions, not just polynomials:

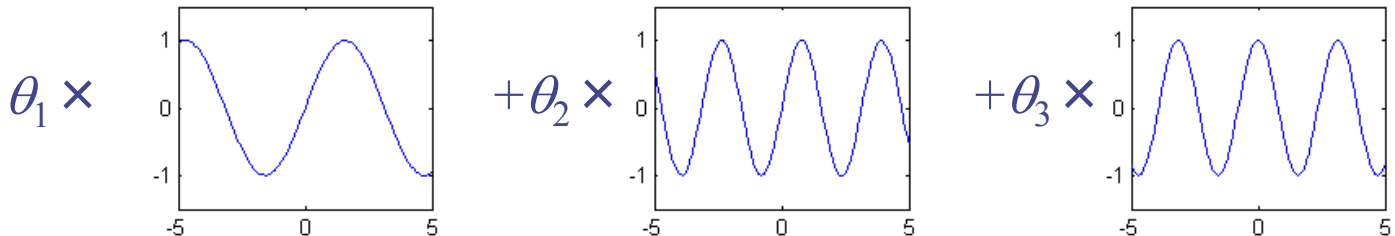
$$f(x; \theta) = \sum_{p=1}^P \theta_p \phi_p(x) + \theta_0$$



- These are generally called **Additive Models**
- Regression adds linear combinations of the basis fn's
- For example: Fourier (sinusoidal) basis

$$\phi_{2k}(x_i) = \sin(kx_i) \quad \phi_{2k+1}(x_i) = \cos(kx_i)$$

- Note, don't have to be a basis per se, usually subset





Patterson
Gimlin

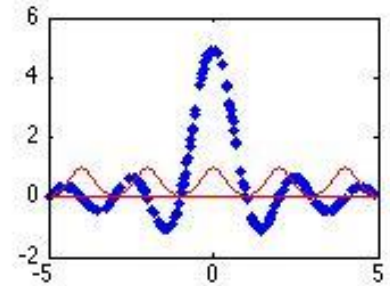
Example: Bigfoot Sightings



Radial Basis Functions (RBF)

- Can act as prototypes of the data itself

$$f(\mathbf{x}; \theta) = \sum_{k=1}^N \theta_k \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x} - \mathbf{x}_k\|^2\right)$$

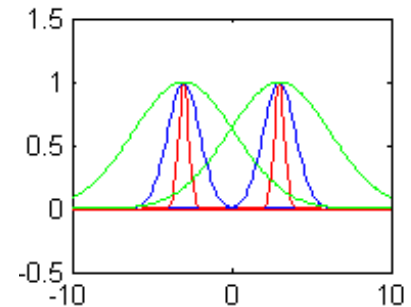
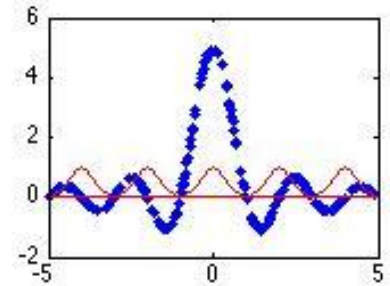


Radial Basis Functions (RBF)

- Can act as prototypes of the data itself

$$f(\mathbf{x}; \theta) = \sum_{k=1}^N \theta_k \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x} - \mathbf{x}_k\|^2\right)$$

- Parameter σ = standard deviation controls how wide bumps are



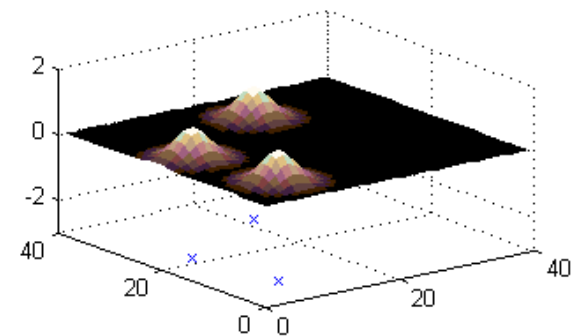
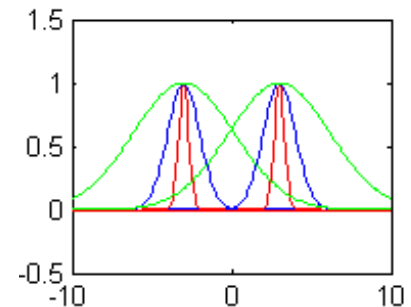
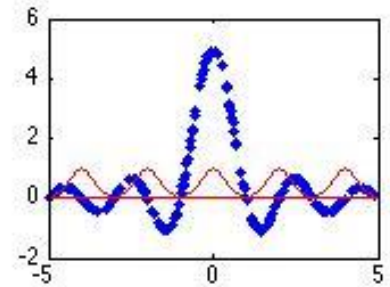
Radial Basis Functions (RBF)

- Can act as prototypes of the data itself

$$f(\mathbf{x}; \theta) = \sum_{k=1}^N \theta_k \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x} - \mathbf{x}_k\|^2\right)$$

- Parameter σ = standard deviation controls how wide bumps are

- Also works in multi-dimensions



Radial Basis Functions

- Training point \rightarrow bump function

$$f(\mathbf{x}; \theta) = \sum_{k=1}^N \theta_k \underbrace{\exp\left(-\frac{1}{\sigma^2} \|\mathbf{x} - \mathbf{x}_k\|^2\right)}_{\phi_k(\mathbf{x})}$$

- Reuse solution from linear regression: $\theta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
- Can view the data instead as \mathbf{X} , a big matrix of size $N \times N$

Radial Basis Functions

- Training point \rightarrow bump function

$$f(\mathbf{x}; \theta) = \sum_{k=1}^N \theta_k \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x} - \mathbf{x}_k\|^2\right) \underbrace{\phi_k(\mathbf{x})}$$

- Reuse solution from linear regression: $\theta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
- Can view the data instead as \mathbf{X} , a big matrix of size $N \times N$

$$\mathbf{X} = \begin{bmatrix} \phi_1(x_1) & \phi_2(x_1) & \cdots & \phi_k(x_1) \\ \phi_1(x_2) & \phi_2(x_2) & \cdots & \phi_k(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_1(x_k) & \phi_2(x_k) & \cdots & \phi_k(x_k) \end{bmatrix}$$

Radial Basis Functions

- Training point \rightarrow bump function

$$f(\mathbf{x}; \theta) = \sum_{k=1}^N \theta_k \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x} - \mathbf{x}_k\|^2\right)$$

- Reuse solution from linear regression: $\theta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
- Can view the data instead as \mathbf{X} , a big matrix of size $N \times N$

$$\mathbf{X} = \begin{bmatrix} \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x}_1 - \mathbf{x}_1\|^2\right) & \cdots & \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x}_1 - \mathbf{x}_k\|^2\right) \\ \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x}_2 - \mathbf{x}_1\|^2\right) & \cdots & \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x}_2 - \mathbf{x}_k\|^2\right) \\ \vdots & \ddots & \vdots \\ \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x}_k - \mathbf{x}_1\|^2\right) & \cdots & \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x}_k - \mathbf{x}_k\|^2\right) \end{bmatrix}$$

Radial Basis Functions

- training point \rightarrow bump function

$$f(\mathbf{x}; \theta) = \sum_{k=1}^N \theta_k \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x} - \mathbf{x}_k\|^2\right)$$

- Reuse solution from linear regression: $\theta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
- Can view the data instead as \mathbf{X} , a big matrix of size $N \times N$

$$\mathbf{X} = \begin{bmatrix} \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x}_1 - \mathbf{x}_1\|^2\right) & \cdots & \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x}_1 - \mathbf{x}_k\|^2\right) \\ \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x}_2 - \mathbf{x}_1\|^2\right) & \cdots & \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x}_2 - \mathbf{x}_k\|^2\right) \\ \vdots & \ddots & \vdots \\ \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x}_k - \mathbf{x}_1\|^2\right) & \cdots & \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x}_k - \mathbf{x}_k\|^2\right) \end{bmatrix}$$

- For RBFs, \mathbf{X} is square and symmetric, so solution is just

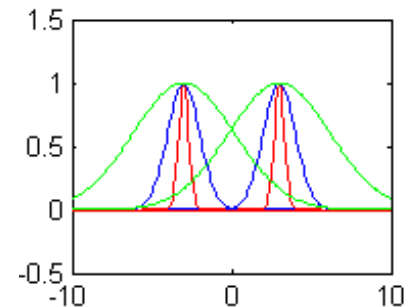
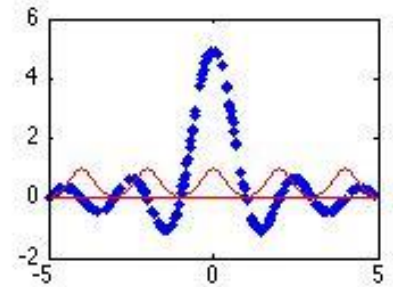
$$\theta^* = \mathbf{X}^{-1} \mathbf{y}$$

Bump Width for RBF

- Can act as prototypes of the data itself

$$f(\mathbf{x}; \theta) = \sum_{k=1}^N \theta_k \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x} - \mathbf{x}_k\|^2\right)$$

- Parameter σ = standard deviation controls how wide bumps are



What happens if too big/small?

How would we know that?

Evaluating Our Learned Function

- We minimized empirical risk to get θ^*
- How well does $f(x; \theta^*)$ perform on future data?
- It should *Generalize* and have low **True Risk**:

$$R_{true}(\theta) = \int P(x, y) \frac{1}{2} (y - \theta^T x)^2 dx dy$$

- Can't compute true risk, instead use **Testing Empirical Risk**
- We randomly split data into training and testing portions

$$\{(x_1, y_1), \dots, (x_N, y_N)\}$$

$$\{(x_{N+1}, y_{N+1}), \dots, (x_{N+M}, y_{N+M})\}$$

- Find θ^* with **training data**:

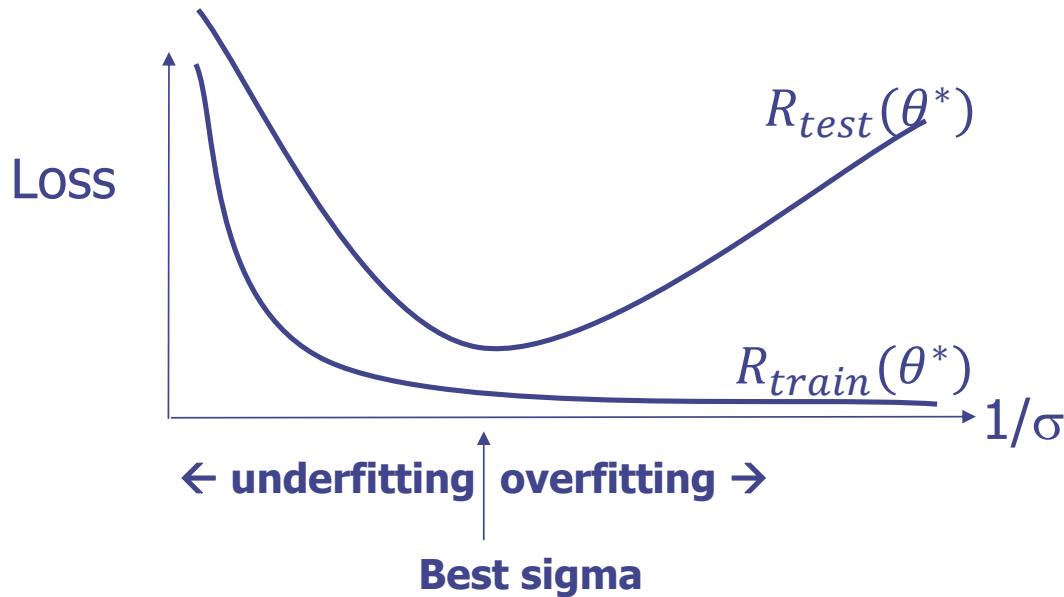
$$R_{train}(\theta) = \frac{1}{2N} \sum_{i=1}^N (y_i - \theta^T x_i)^2$$

- Evaluate it with **testing data**:

$$R_{train}(\theta) = \frac{1}{2M} \sum_{i=N+1}^{N+M} (y_i - \theta^T x_i)^2$$

Crossvalidation

- Try fitting with different sigma radial basis function widths
- Select sigma which gives lowest $R_{test}(\theta^*)$



- Think of sigma as a measure of the simplicity of the model
- Thinner RBFs are more flexible and complex