# CPSC 340:
# Machine Learning and Data Mining

Regularization

# Admin
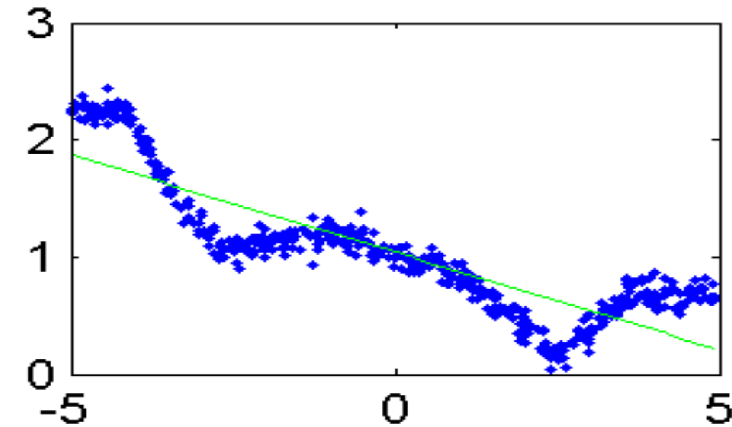
- <span style="color:red">Assignment 2</span>:
  - 1 late day to hand it in today, 2 for tomorrow, 3 for Wednesday.
- <span style="color:red">Assignment 3</span> is out.
  - Due next Friday (right before the break)
- Tutorials this week: assignment 3 practice
- Real-time feedback system is back up

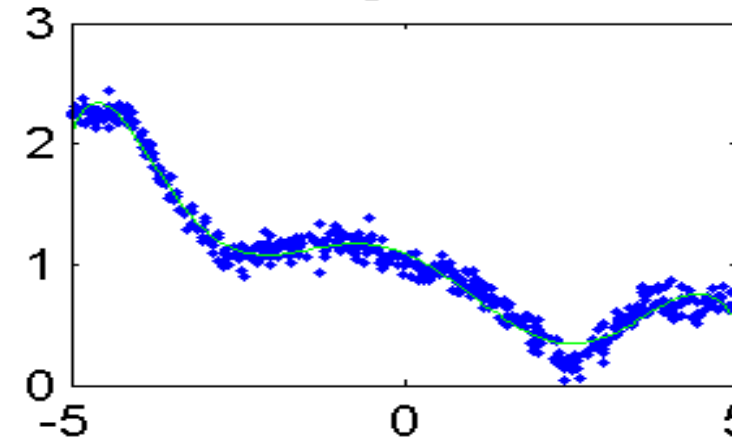# Last Time: Normal Equations and Change of Basis

- Last time we derived normal equations:

$$X^T X w = X^T y$$

  – Solutions 'w' minimize squared error in linear model.

- We also discussed change of basis:
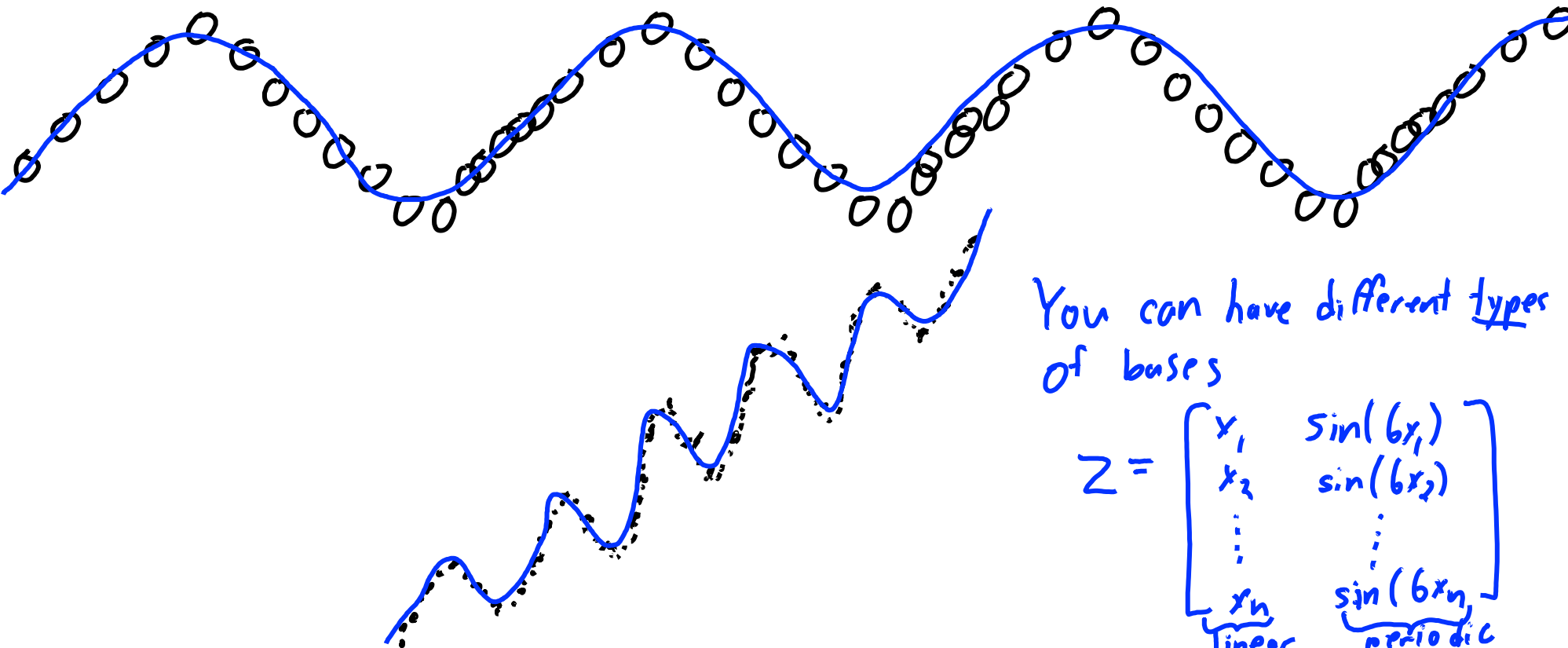
  – E.g., polynomial basis:

$$\text{Replace } X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \text{ with } Z = \begin{bmatrix} 1 & x_1 & (x_1)^2 & \cdots & (x_1)^p \\ 1 & x_2 & (x_2)^2 & \cdots & (x_2)^p \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & (x_n)^2 & \cdots & (x_n)^p \end{bmatrix}$$

  – Let's you fit non-linear models with linear regression.

$$y_i = w^T z_i = w_0 + w_1 x_i + w_2 x_i^2 + w_3 x_i^3 + \cdots + w_p x_i^p$$

Degree 7

3

# Parametric vs. Non-Parametric Bases

- Polynomials are not the only possible bases:
  - Exponentials, logarithms, trigonometric functions, etc.
  - The right basis will vastly improve performance.

For periodic data we might use

$$Z = \begin{bmatrix} \sin(x_1) \\ \sin(x_2) \\ \vdots \\ \sin(x_n) \end{bmatrix}$$

You can have different types of bases

$$Z = \begin{bmatrix} x_1 & \sin(6x_1) \\ x_2 & \sin(6x_2) \\ \vdots & \vdots \\ x_n & \sin(6x_n) \end{bmatrix}$$

linear    periodic

# Parametric vs. Non-Parametric Bases

- Polynomials are not the only possible bases:
  - Exponentials, logarithms, trigonometric functions, etc.
  - The right basis will vastly improve performance.
  - But the right basis may not be obvious.
- What happens if we use the wrong basis?
  - As 'n' increases, we can fit 'w' more accurately.
  - But eventually more data doesn't help if basis isn't "flexible" enough.
- Alternative is non-parametric bases:
  - Size of basis (number of features) grows with 'n'.
  - Model gets more complicated as you get more data.
  - You can model very complicated functions where you don't know the right basis.

# Non-Parametric Basis: RBFs

- **Radial basis functions** (RBFs):
  - Non-parametric bases that depend on distances to training points.

$$\text{Replace} \quad X = \begin{bmatrix} \quad \\ \quad \\ \quad \\ \quad \end{bmatrix} \Big\} n \qquad \text{by} \qquad Z = \begin{bmatrix} g(\|x_1 - x_1\|) & g(\|x_1 - x_2\|) & \cdots & g(\|x_1 - x_n\|) \\ g(\|x_2 - x_1\|) & g(\|x_2 - x_2\|) & \cdots & g(\|x_2 - x_n\|) \\ \vdots & & \ddots & \vdots \\ g(\|x_n - x_1\|) & g(\|x_n - x_2\|) & \cdots & g(\|x_n - x_n\|) \end{bmatrix} \Big\} n$$

(under $X$: $\underbrace{\quad}_{d}$)   (under $Z$: $\underbrace{\quad}_{n}$)

  - Most common 'g' is Gaussian RBF:

$$g(\alpha) = \exp\left(-\frac{\alpha^2}{2\sigma^2}\right)$$

→ Parameter

$g(\|x_1 - x_i\|)$

$x_i \qquad x_1$

- Variance $\sigma^2$ controls influence of nearby points.
- This affects fundamental trade-off (set it using a validation set).

Do we need $\sigma\sqrt{2\pi}$?
- No because $w^\top x_i = \left(\frac{1}{\beta}w\right)^\top(\beta x_i)$

# Non-Parametric Basis: RBFs

- Radial basis functions (RBFs):
  - Non-parametric bases that depend on distances to training points.

$$\text{Replace} \quad X = \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix} \Big\} n \quad \text{by} \quad Z = \begin{bmatrix} g(\|x_1 - x_1\|) & g(\|x_1 - x_2\|) & \cdots & g(\|x_1 - x_n\|) \\ g(\|x_2 - x_1\|) & g(\|x_2 - x_2\|) & \cdots & g(\|x_2 - x_n\|) \\ \vdots & & & \vdots \\ g(\|x_n - x_1\|) & g(\|x_n - x_2\|) & \cdots & g(\|x_n - x_n\|) \end{bmatrix} \Big\} n$$

$$\underbrace{\phantom{xxxxxx}}_{d} \qquad \underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxxxx}}_{n}$$

$$\text{To make predictions on} \quad \hat{X} = \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix} \Big\} t \quad \text{use} \quad \hat{Z} = \begin{bmatrix} \\ g(\|\hat{x}_i - x_j\|) \\ \\ \end{bmatrix} \Big\} t$$

$$\underbrace{\phantom{xxxx}}_{d} \qquad\qquad \underbrace{\phantom{xxxxxxxxxxxxxx}}_{n}$$

I.e.,
$$y_i = w^T \hat{z}_i$$
or $\hat{y} = Zw$

Number of "features" is number of training examples.

# Non-Parametric Basis: RBFs

Cubic basis: $y_i = w_0$  $+ w_1$  $+ w_2$  $+ w_3$  $+ w_4$ 

Polynomial basis represents function as sum of global polynomials.

Gaussian RBFs: $y_i = w_0$  $+ w_1$  $+ w_2$  $+ w_3$  $+ w_4$ 

Gaussian RBFs represent function as sum of local "bumps"

- Gaussian RBFs are universal approximators (compact subets of $\mathbb{R}^d$)
  - Can approximate any continuous function to arbitrary precision.
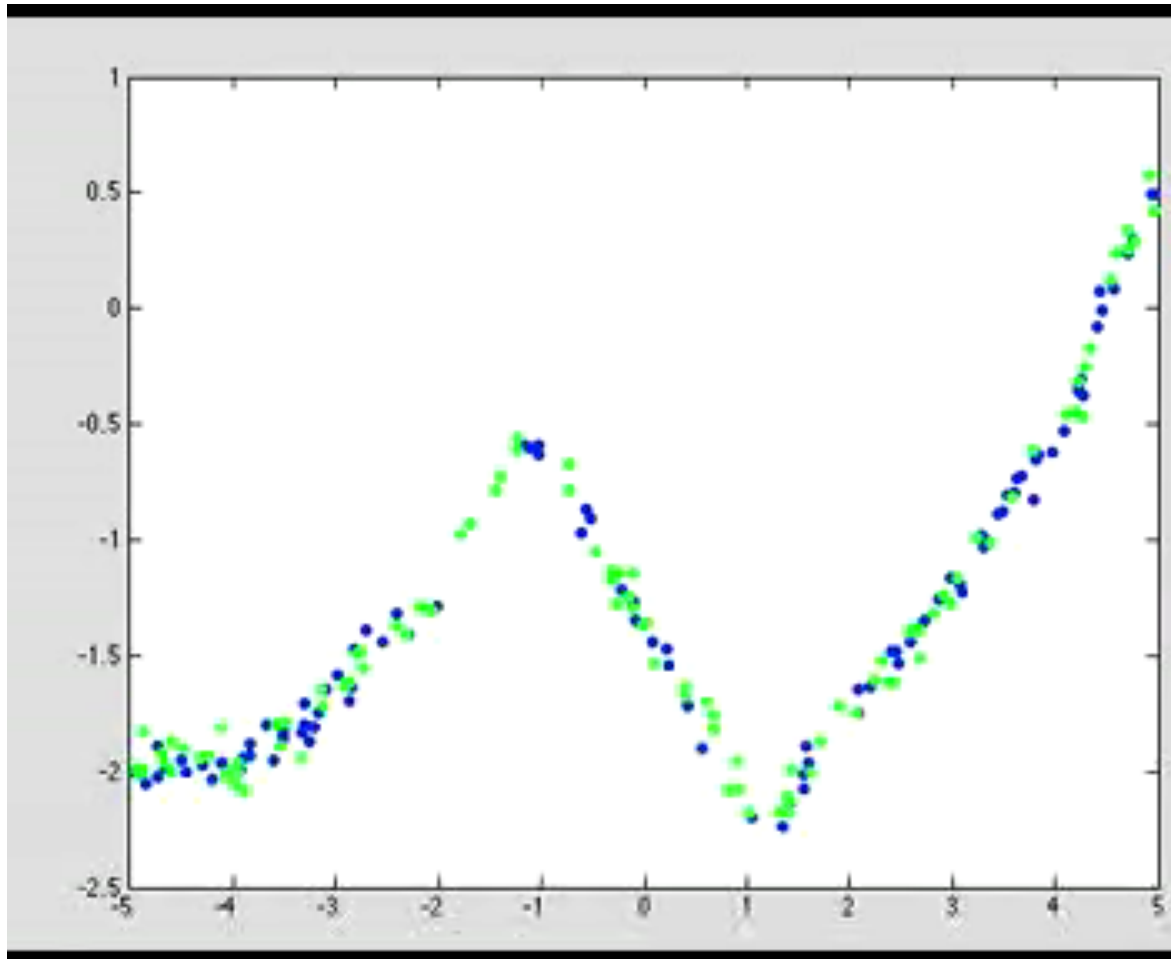  - Achieve irreducible error as 'n' goes to infinity.

# Interpolation vs. Extrapolation



polynomial basis becomes polynomial away from data

Gaussian RBFs go to _zero_ away from data.

9

# Non-Parametric Basis: RBFs

- Least squares with Gaussian RBFs for different σ values:



Could add __bias__ and linear basis:

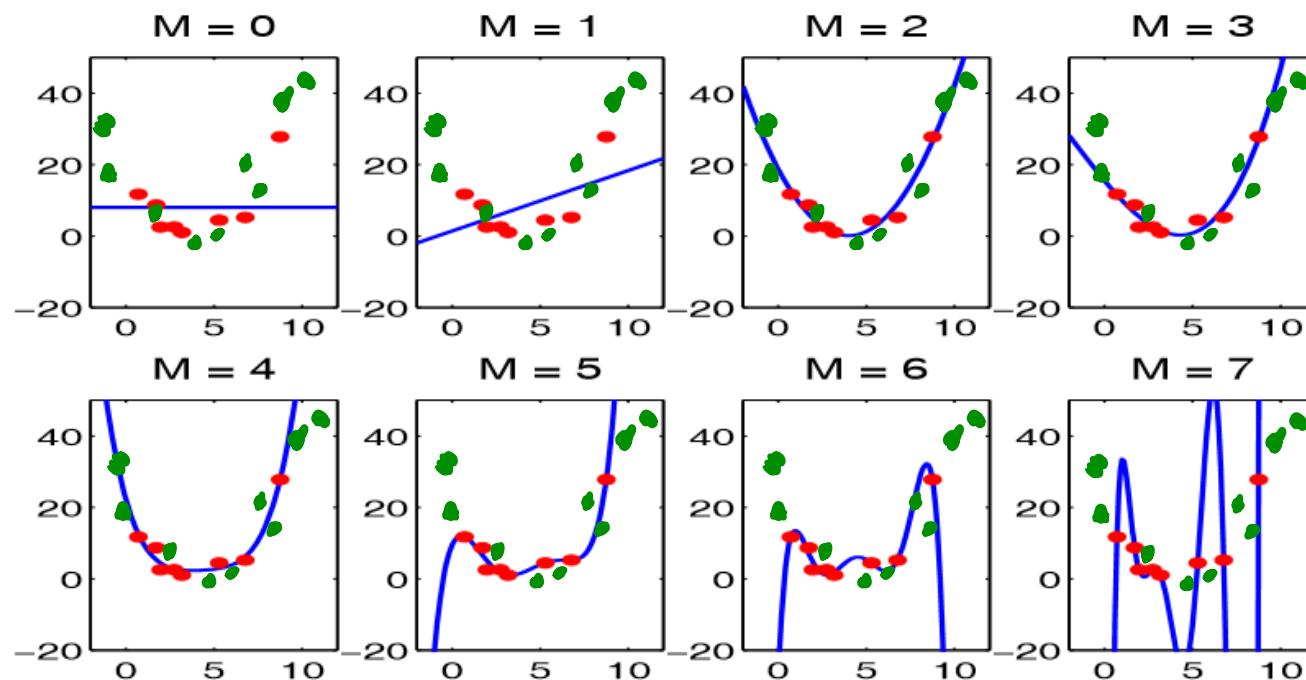$$Z = \begin{bmatrix} 1 & -x_1- & g(\|x_1-x_1\|) & \cdots & g(\|x_1-x_n\|) \\ 1 & -x_2- & & & \vdots \\ 1 & -x_3- & \vdots & & \\ \vdots & \vdots & & & \vdots \\ 1 & -x_n- & g(\|x_1-x_n\|) & \cdots & g(\|x_n-x_n\|) \end{bmatrix}$$

$\underbrace{\phantom{1}}_{1}$ $\underbrace{\phantom{-x_n-}}_{d}$ $\underbrace{\phantom{g(\|x_1-x_n\|)\cdots g(\|x_n-x_n\|)}}_{n}$

This reverts to <u>linear regression</u> instead of 0 away from data.

# Last Time: Polynomial Degree and Training vs. Testing

- As the polynomial degree increases, the training error goes down.
- But training error becomes worse approximation test error.



- Same effect as we decrease variance in Gaussian RBF.
- But what if we need a complicated model?

# Controlling Complexity

- Usually <span style="color:red">"true" mapping from $x_i$ to $y_i$ is complex.</span>
  - Might need high-degree polynomial or small $\sigma^2$ in RBFs.
- But <span style="color:red">complex models can overfit</span>.
- So what do we do???

- There are many possible answers:
  - <span style="color:blue">Model averaging</span>: average over multiple models to decrease variance.
  - <span style="color:blue">Regularization</span>: add a <span style="color:green">penalty on the complexity</span> of the model.

# L2-Regularization

- Standard regularization strategy is L2-regularization:

$$f(w) = \frac{1}{2}\|Xw - y\|^2 + \frac{\lambda}{2}\|w\|^2 \qquad or \qquad f(w) = \frac{1}{2}\sum_{i=1}^{n}(w^T x_i - y_i)^2 + \frac{\lambda}{2}\sum_{j=1}^{d} w_j^2$$

"lambda"

- Intuition: large $w_j$ tend to lead to overfitting (cancel each other).

- So minimize squared error plus penalty on L2-norm of 'w'.

  – This objective balances getting low error vs. having small slope 'w'.

    - Training error will increase *because you're no longer minimizing it*.

    - But reduces overfitting.

  – Regularization parameter λ > 0 controls "strength" of regularization.

    - Large λ puts large penalty on slope.

    - There is such a thing as "too much" regularization: $\lambda \rightarrow \infty$ will cause w=0

13

# L2-Regularization

- Standard regularization strategy is L2-regularization:

$$f(w) = \frac{1}{2}\|Xw - y\|^2 + \frac{\lambda}{2}\|w\|^2 \qquad or \qquad f(w) = \frac{1}{2}\sum_{i=1}^{n}(w^T x_i - y_i)^2 + \frac{\lambda}{2}\sum_{j=1}^{d} w_j^2$$
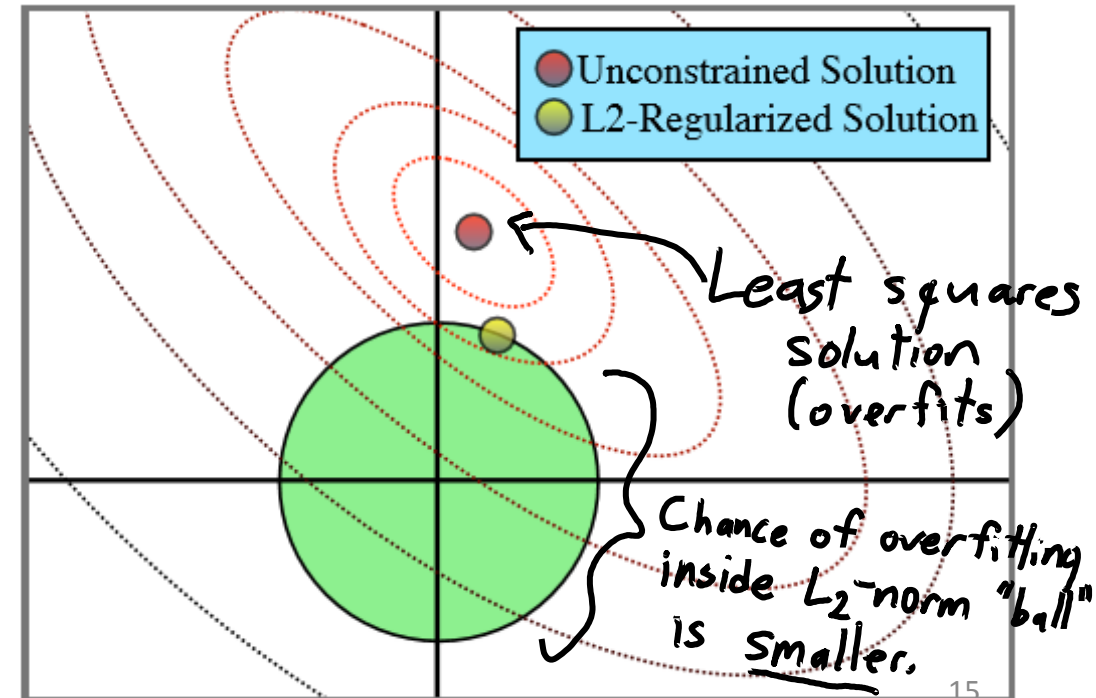
- In terms of fundamental trade-off:
  – Regularization increases training error.
  – Regularization makes training error better approximation of test error.
- How should you choose λ?
  – Theory: as 'n' grows λ should be in the range $O(1)$ to $O(n^{1/2})$.
  – Practice: optimize validation set or cross-validation error.
    - This almost always decreases the test error.

# L2-Regularization

- Standard regularization strategy is L2-regularization:

$$f(w) = \frac{1}{2}\|Xw - y\|^2 + \frac{\lambda}{2}\|w\|^2 \qquad \text{or} \qquad f(w) = \frac{1}{2}\sum_{i=1}^{n}(w^T x_i - y_i)^2 + \frac{\lambda}{2}\sum_{j=1}^{d} w_j^2$$

- Equivalent to minimizing squared error with L2-norm constraint:

- Connection to Occam's razor
  - Small values in 'w' are "simpler" models
  - L2-regularization favors small 'w'

- Regularization is a way of incorporating prior knowledge



Unconstrained Solution
L2-Regularized Solution

Least squares solution (overfits)

Chance of overfitting inside L2-norm "ball" is smaller.

15

# Why use L2-Regularization?

- Almost always decreases test error
- Intuition: try to make the objective function reflect test error
  - The original objective function was just *training error*
  - Create an optimization problem that you actually want to solve
- But here are 6 more reasons (for linear regression):
  1. Solution 'w' is unique.
  2. $X^TX$ does not need to be invertible.
  3. Less sensitive to changes in X or y (related to uniqueness).
  4. Makes algorithms for computing 'w' converge faster.
  5. Stein's paradox: if d ≥ 3, 'shrinking' moves us closer to 'true' w on average.
     - In other words, it's a good idea even if the prior knowledge is wrong (!!) – see bonus slides
  6. Worst case: just set λ small and get the same performance.

# RBFs, Regularization, and Validation

- A model that is hard to beat:
  - RBF basis with L2-regularization and cross-validation to choose $\sigma$ and $\lambda$.
  - Flexible non-parametric basis, magic of regularization, and tuning for test error!

Example:
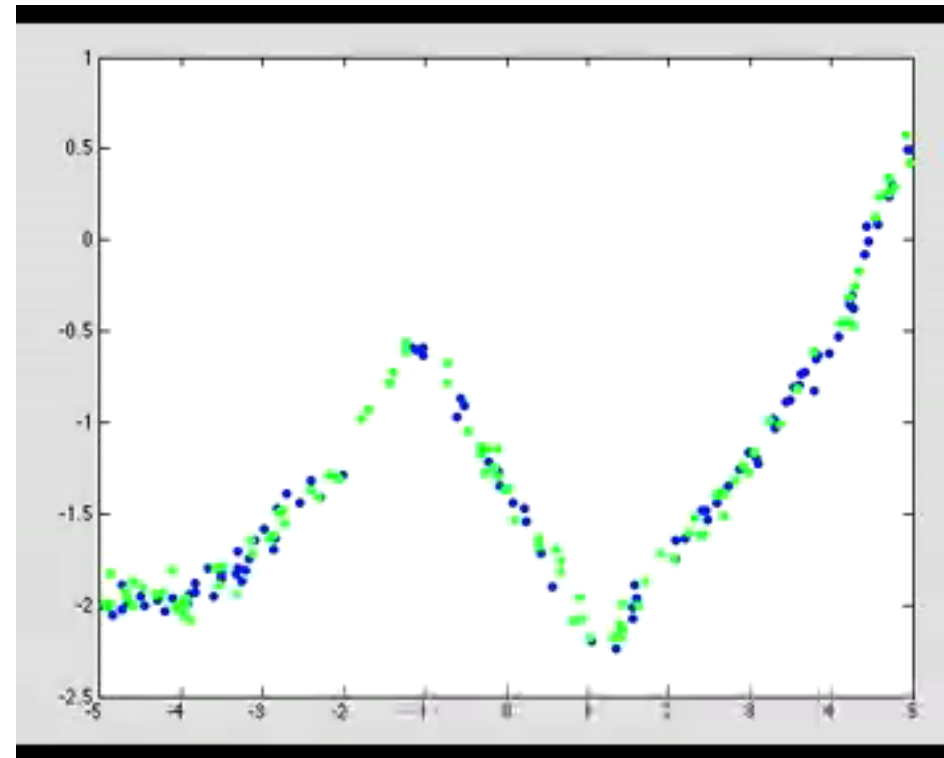Find regularized value of 'w' for particular $\lambda$ and $\sigma$ by

minimizing $f(w) = \frac{1}{2} \| Zw - y \|^2 + \frac{\lambda}{2} \| w \|^2$

$\hookleftarrow$ RBF basis. with variance $\sigma$

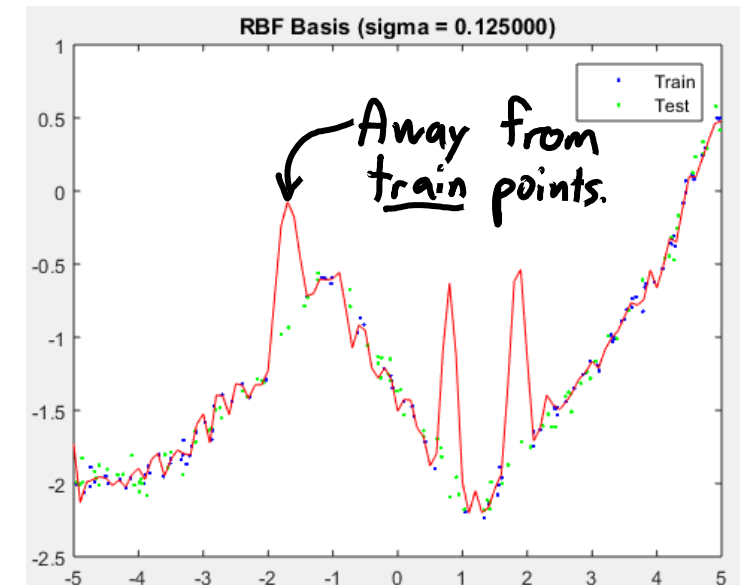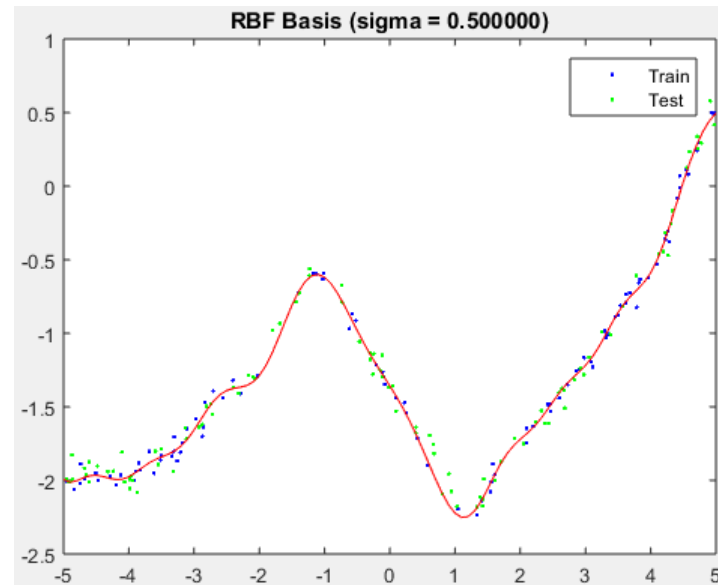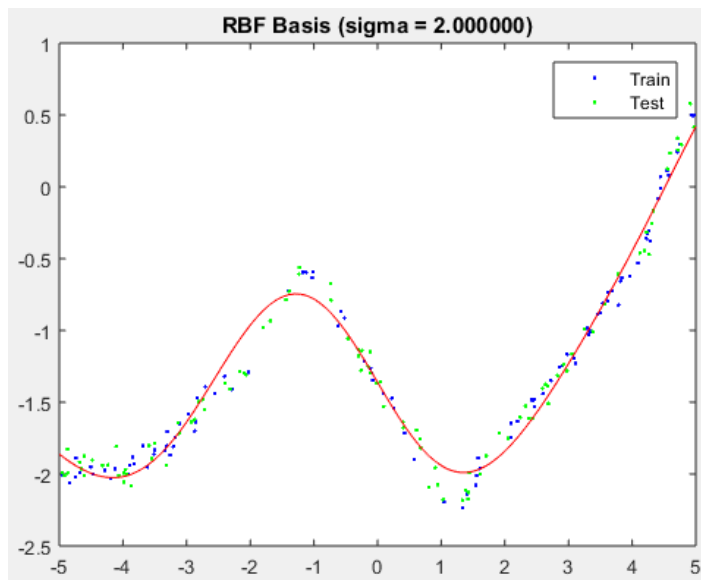And choose $\lambda$ and $\sigma$ to minimize

$\frac{1}{2} \| \hat{Z} \hat{w} - \hat{y} \|^2$

Validation set

$\longrightarrow$ Regularized value of w

# RBFs, Regularization, and Validation

- A model that is hard to beat:
  - RBF basis with L2-regularization and cross-validation to choose $\sigma$ and $\lambda$.
  - Flexible non-parametric basis, magic of regularization, and tuning for test error!



  - Can add bias or linear/poly basis to do better away from data.
  - Expensive at test time: need distance to all training examples.

# Summary

- **Radial basis functions:**
  - Non-parametric bases that can model any function.
- **Regularization:**
  - Adding a penalty on model complexity.
  - Improves test error because it is magic.
- **L2-regularization:** penalty on L2-norm of regression weights 'w'.

- Next time:
  - Going downhill…