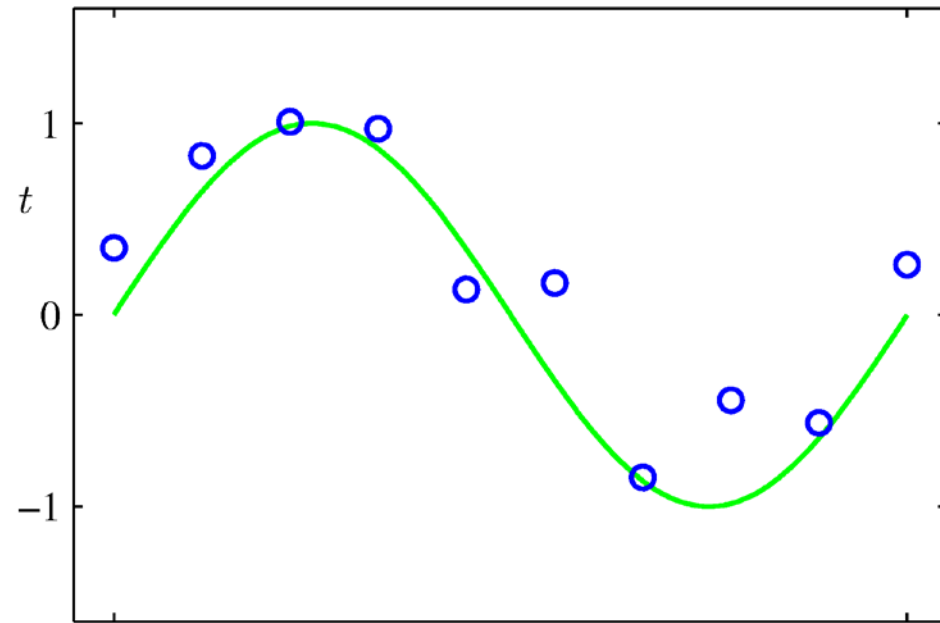# Notes on regularization

CS434
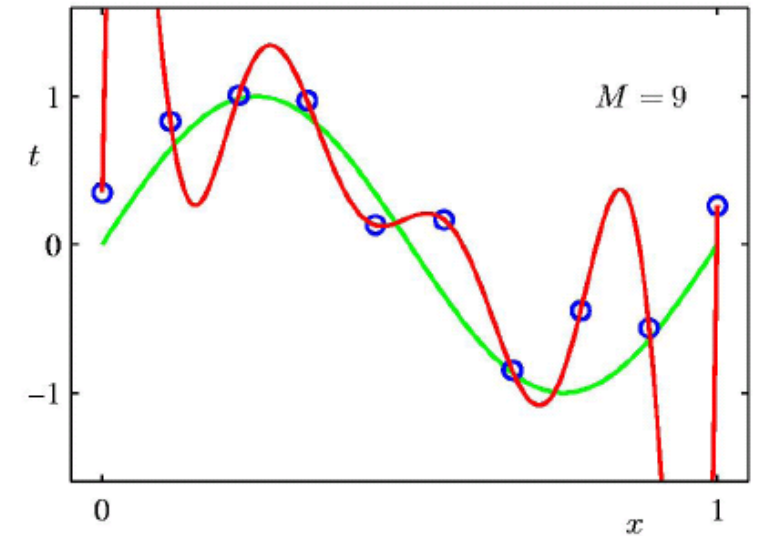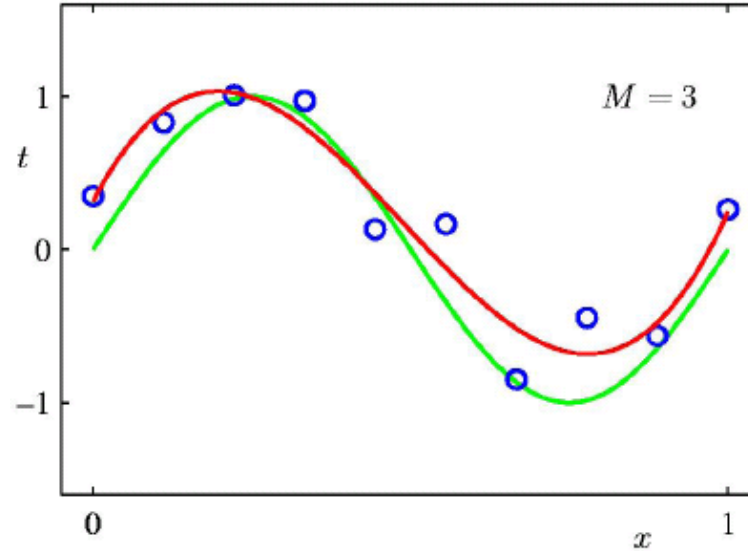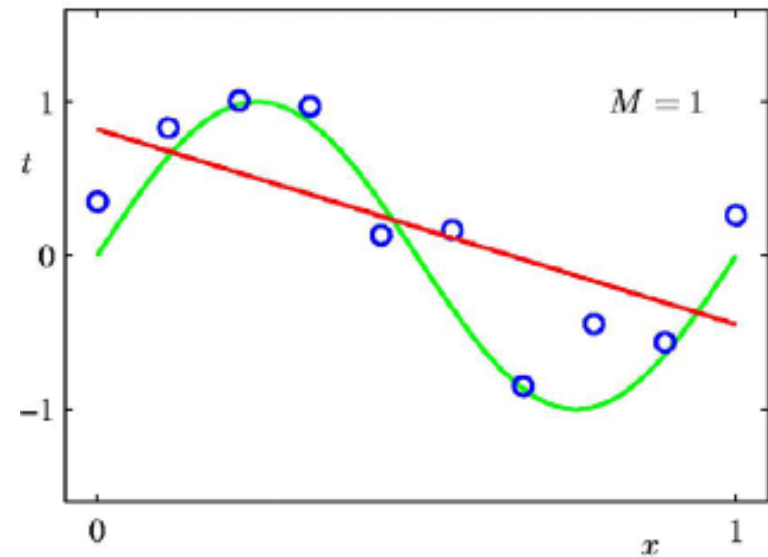
# A regression example: Polynomial Curve Fitting



$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \ldots + w_M x^M = \sum_{j=0}^{M} w_j x^j$$

- In this example, there is only one feature $x$. We learn a function of M-order polynomial
- Alternatively, we could also view this as linear regression using $(1, x, x^2, \ldots, x^M)$ as the features.
- Note that this new feature space is derived from the original input $x$
- Such derived features are often referred to as the basis functions

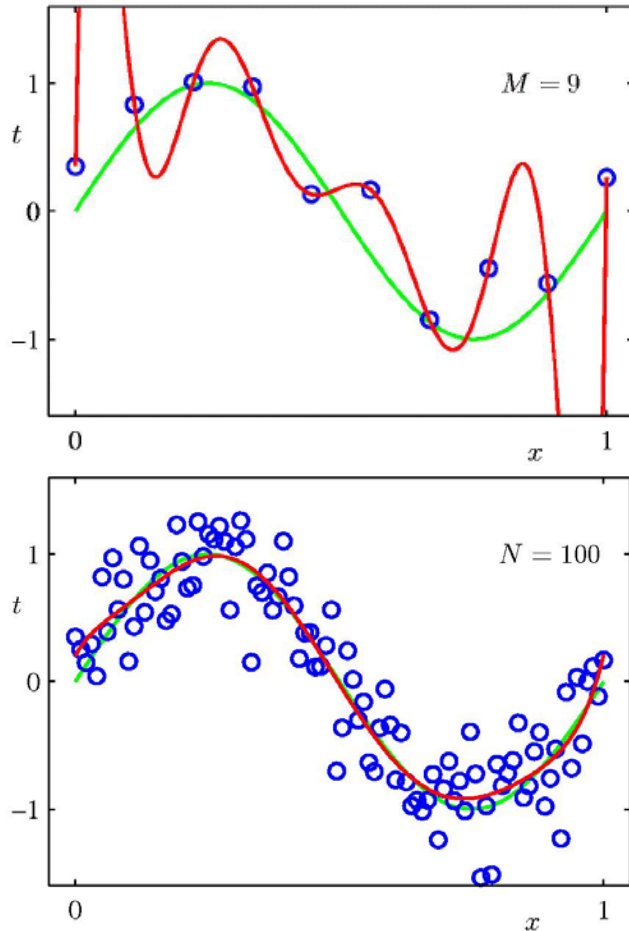# Consider different choices for M



$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \ldots + w_M x^M = \sum_{j=0}^{M} w_j x^j$$

- Larger M leads to higher model complexity
- Given 10 data points, if M=9, we can fit the training data perfectly – severely overfitting

# Over-fitting issue



- What can we do to curb over-fitting
  - Use less complex model
  - Use more training examples
  - **Regularization**

In linear regression, overfitting can often be characterized by large weights

| | M = 0 | M = 1 | M = 3 | M = 9 |
|---|---|---|---|---|
| $w_0$ | 0.19 | 0.82 | 0.31 | 0.35 |
| $w_1$ | | -1.27 | 7.99 | 232.37 |
| $w_2$ | | | -25.43 | -5321.83 |
| $w_3$ | | | 17.37 | 48568.31 |
| $w_4$ | | | | -231639.30 |
| $w_5$ | | | | 640042.26 |
| $w_6$ | | | | -1061800.52 |
| $w_7$ | | | | 1042400.18 |
| $w_8$ | | | | -557682.99 |
| $w_9$ | | | | 125201.43 |

# Regularized Linear Regression

- Consider the following loss function:

$$E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$$

Data term + Regularization term (penalize complex models)

$$\sum_{i=1}^{N}(y_i - \mathbf{w}^T\mathbf{x}_i)^2 + \lambda \sum_{j=0}^{M}|w_j|^q$$

Encourage small weight values

| | M = 0 | M = 1 | M = 3 | M = 9 |
|---|---|---|---|---|
| $w_0$ | 0.19 | 0.82 | 0.31 | 0.35 |
| $w_1$ | | -1.27 | 7.99 | 232.37 |
| $w_2$ | | | -25.43 | -5321.83 |
| $w_3$ | | | 17.37 | 48568.31 |
| $w_4$ | | | | -231639.30 |
| $w_5$ | | | | 640042.26 |
| $w_6$ | | | | -1061800.52 |
| $w_7$ | | | | 1042400.18 |

# L2 Regularized Linear Regression

- With the SSE loss and a **quadratic regularizer**, we get

$$\frac{1}{2}\sum_{i=1}^{N}(y_i - \mathbf{w}^T\mathbf{x}_i)^2 + \frac{\lambda}{2}\mathbf{w}^T\mathbf{w}$$

- which is minimized by

$$\mathbf{w} = (\lambda\boldsymbol{I} + \mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^TY$$

- $\lambda$: regularization coefficient, which controls the trade-off between model complexity and the fit to the data
  - Larger $\lambda$ encourages simple model (driving more elements of $\mathbf{w}$ to 0)
  - Small $\lambda$ encourages better fit of the data (driving SSE to zero)
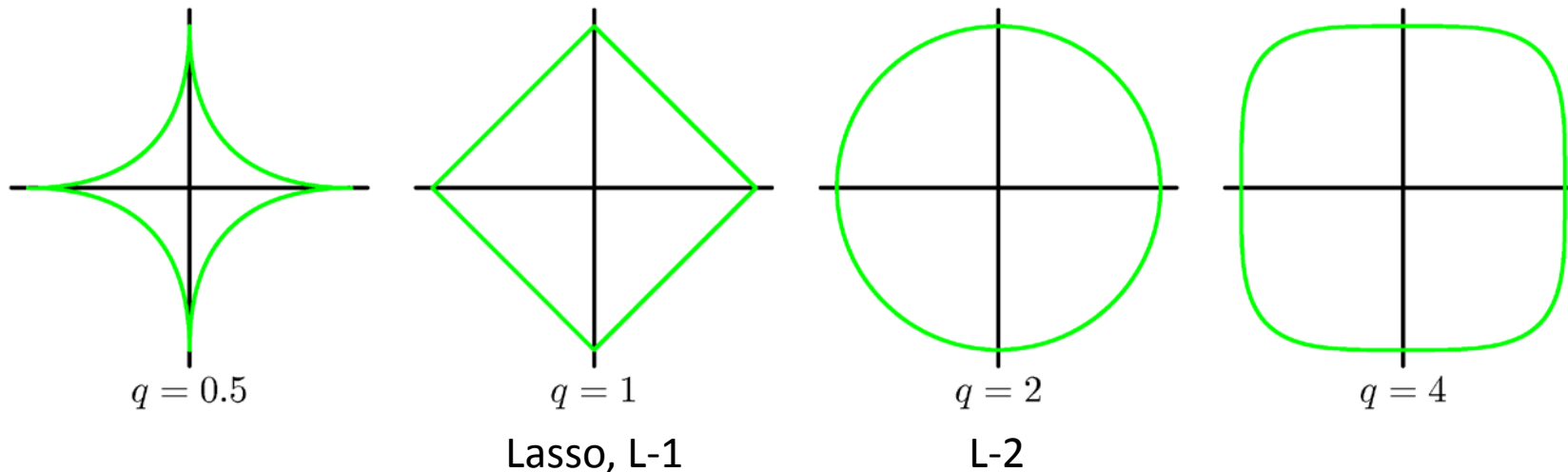
# More Regularizations

$$\sum_{i=1}^{N} (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \sum_{j=0}^{M} |w_j|^q$$

Equivalent to minimizing SSE subject to $\sum_{i=0}^{M} |w_i|^q \leq \epsilon$
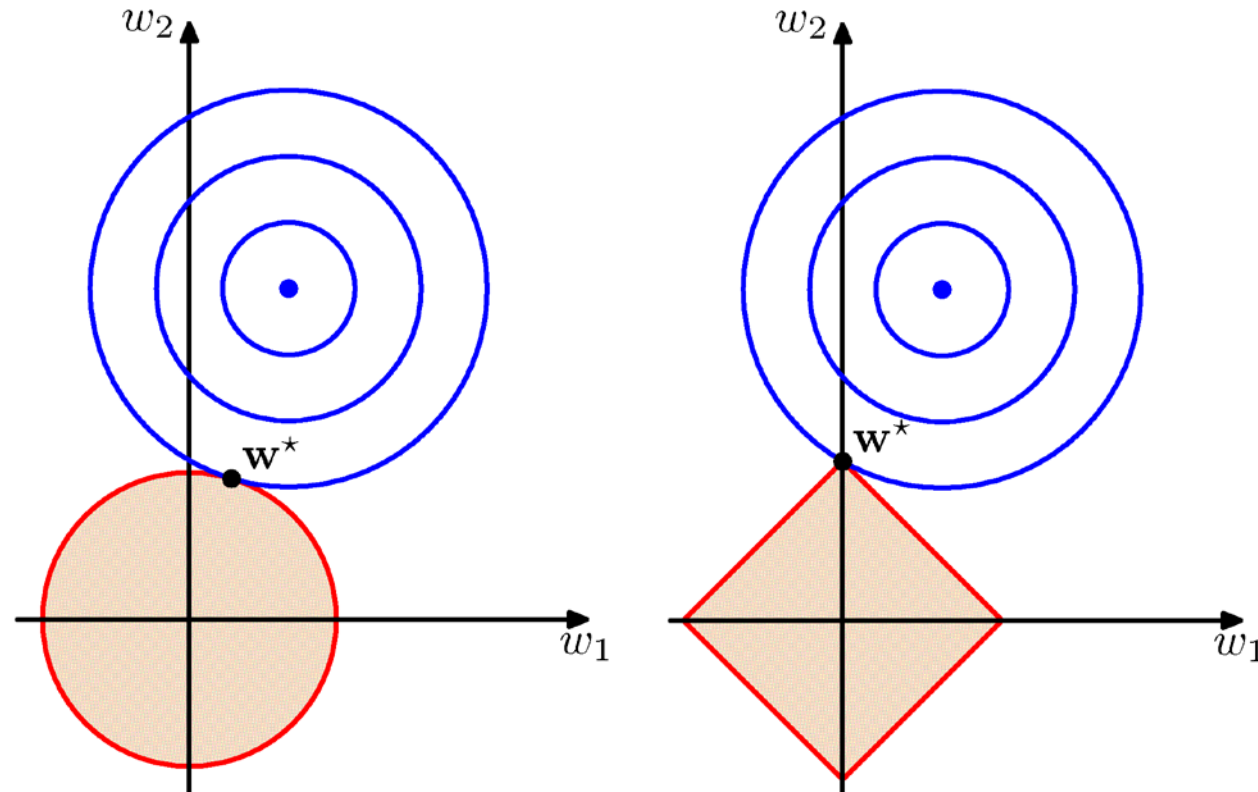
A good explanation of this equivalence is provided here:
http://math.stackexchange.com/questions/335306/why-are-additional-constraint-and-penalty-term-equivalent-in-ridge-regression



$q = 0.5$    $q = 1$    $q = 2$    $q = 4$

Lasso, L-1    L-2

Shape is determined by $q$, size determined by $\lambda$

# Regularized Linear Regression

- Lasso ($q = 1$) tends to generate sparser solutions (majority of the weights shrink to zero) than a quadratic regularizer ($q = 2$, often called ridge regression).

# Commonly used regularizers

- L-2 regularization $\qquad \sum_{i=1}^{N}(y_i - \mathbf{w}^T\mathbf{x}_i)^2 + \lambda\sum_{j=0}^{M}w_j{}^2$

  Poly-time close-form solution
  Curbs overfitting but does not produce sparse solution

- L-1 regularization $\qquad \sum_{i=1}^{N}(y_i - \mathbf{w}^T\mathbf{x}_i)^2 + \sum_{j=0}^{M}|w_j|$

  Poly-time approximation algorithm
  Sparse solution – potentially many zeros in $\mathbf{w}$

- L-0 regularization $\qquad \sum_{i=1}^{N}(y_i - \mathbf{w}^T\mathbf{x}_i)^2 + \sum_{j=0}^{M}I(w_j \neq 0)$

  Seek to identify optimal feature subset
  NP-complete problem!