
Regression, function approximation

Some Common Themes

Illustrated through (Generalized) MLR

Readings/Notation: I'll closely follow Bishop Ch 3.1, 3.2, which uses machine learning notation: parameters are w 's (for weights), dependent variable is " t " for target, and model produces output " y ".

Function Approximation / Regression/Prediction

- A predictive modeling technique
 - Given:
 - A set of input (AI) /independent (math)/ explanatory or predictor (stats) variables X
 - corresponding (set of) output/dependent/response variables Y
 - Build: a model relating X to Y
 - single value for Y given X (more common)
 - e.g. $E[Y|X]$, the “regression of y on X .”
 - Assumes $Y = \text{function of } X + (\text{zero-mean, symmetric}) \text{ noise}$
 - Add Confidence Interval
 - Distribution of Y given X

Parametric Models

Determine **functional form** of model (e.g. polynomials)

- “learn” the parameters (weights) of the model using the training data.

- **Example:** linear regression

- **Generalize:** linear combination of basis functions (basis function expansion)

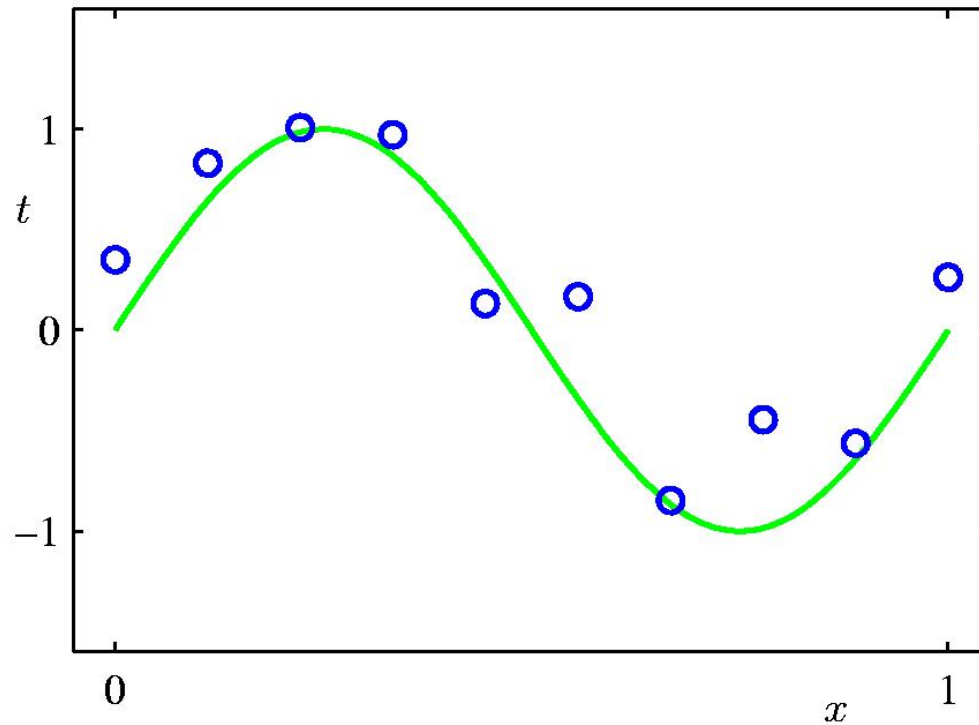
$$y(x, \mathbf{w}) = \sum_{i=0}^M w_i \phi_i(x) = \mathbf{w}^T \boldsymbol{\phi}(x)$$

- Special Case: linear regression.
- Special Case: polynomial: (with scalar x)

$$y(x, \mathbf{w}) = w_0 + w_1 x + \dots + w_M x^M$$

so that the basis functions are given by $\phi_i(x) = x^i$

Polynomial Curve Fitting



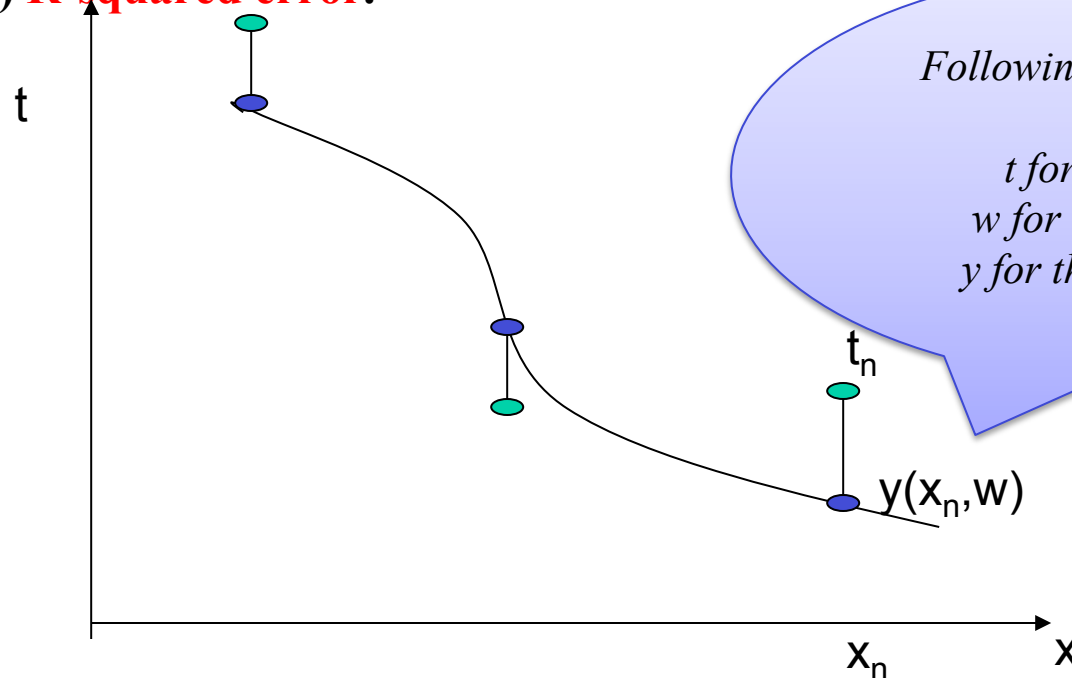
$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

Least Squares

- Minimize sum-of-squares **error (SSE)** (**t's are the target values**)

$$E(\mathbf{w}) = \sum_{n=1}^N \{\mathbf{w}^T \phi(x_n) - t_n\}^2$$

Best interpretation? Consider Root Mean Squared Error (RMSE) or (adjusted) **R-squared error.**



*Following Bishop (06), am
using
t for target values
w for the parameters
y for the model output*

Least Squares Solution*

- Exact closed-form minimizer (ML solution)

$$\mathbf{w}^* = (\Phi^T \Phi)^{-1} \Phi^T \vec{t}$$

where $\vec{t} = (t_1, \dots, t_N)^T$

– “Pseudo-inverse solution”

and Φ is the *design matrix* given by

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \cdots & \phi_M(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \cdots & \phi_M(\mathbf{x}_2) \\ \vdots & \vdots & \vdots \\ \phi_0(\mathbf{x}_N) & \cdots & \phi_M(\mathbf{x}_N) \end{pmatrix}$$

Takeaway: direction solution involves inversion of a matrix

Explicitly shows collinearity problem

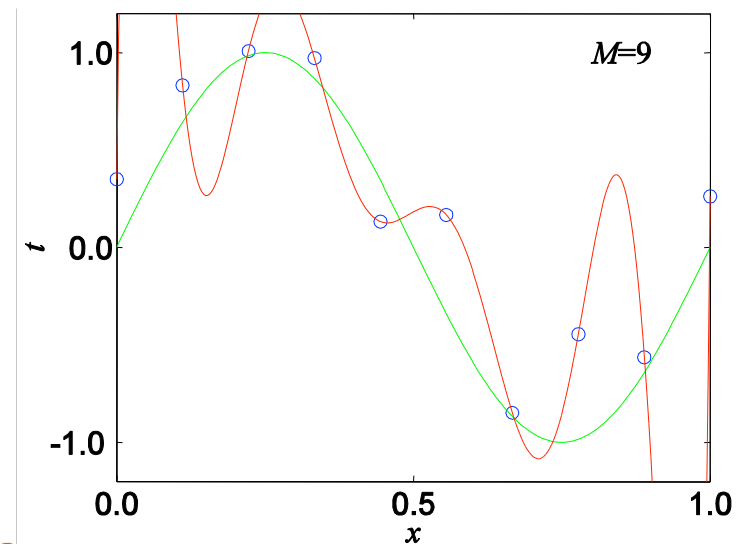
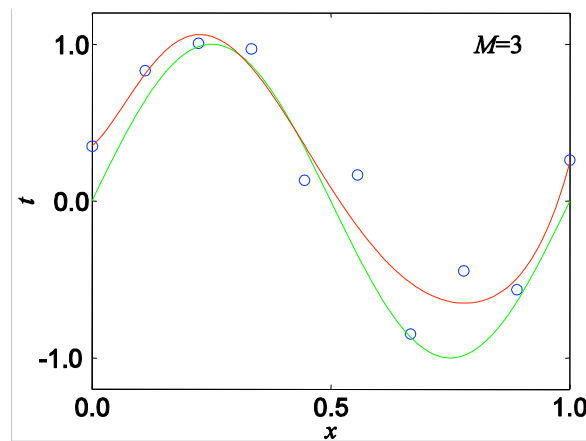
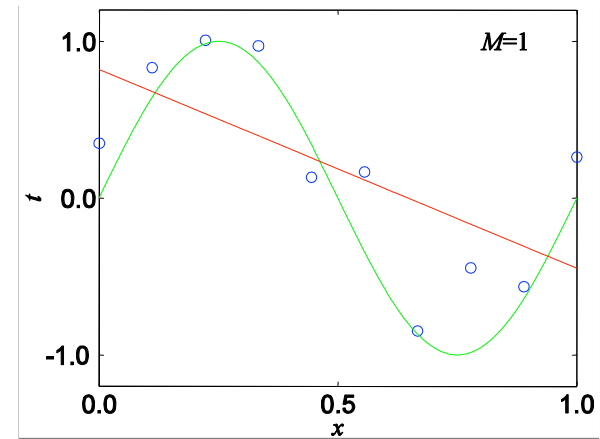
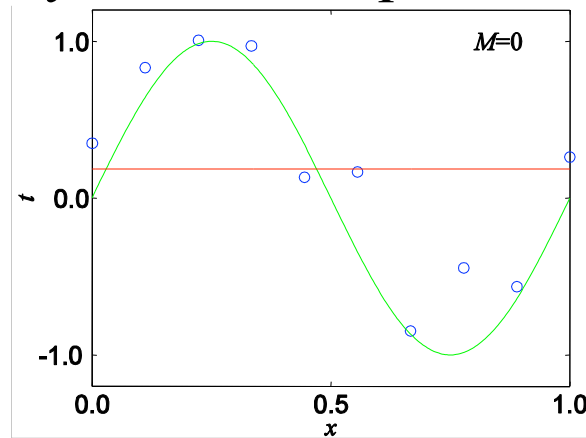
Multiple outputs?

$$\mathbf{w}_k^* = (\Phi^T \Phi)^{-1} \Phi^T \vec{t}_k$$

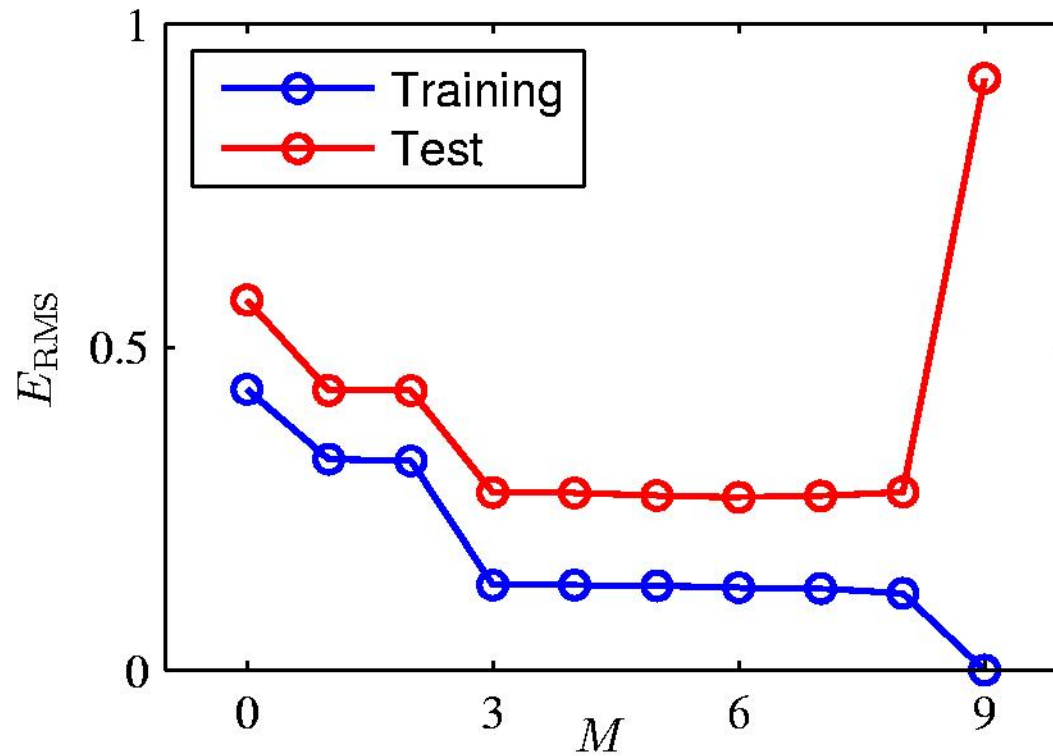
(pseudo-inverse portioned part shared by all outputs; rest de-coupled.)

Model Complexity and Overfitting

- “Noisy sine” example from Chris Bishop



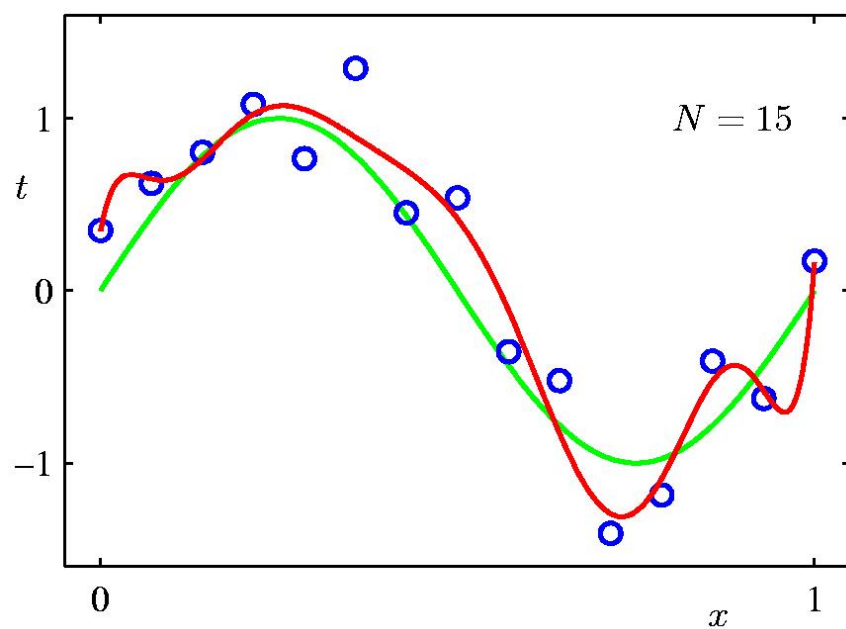
Over-fitting



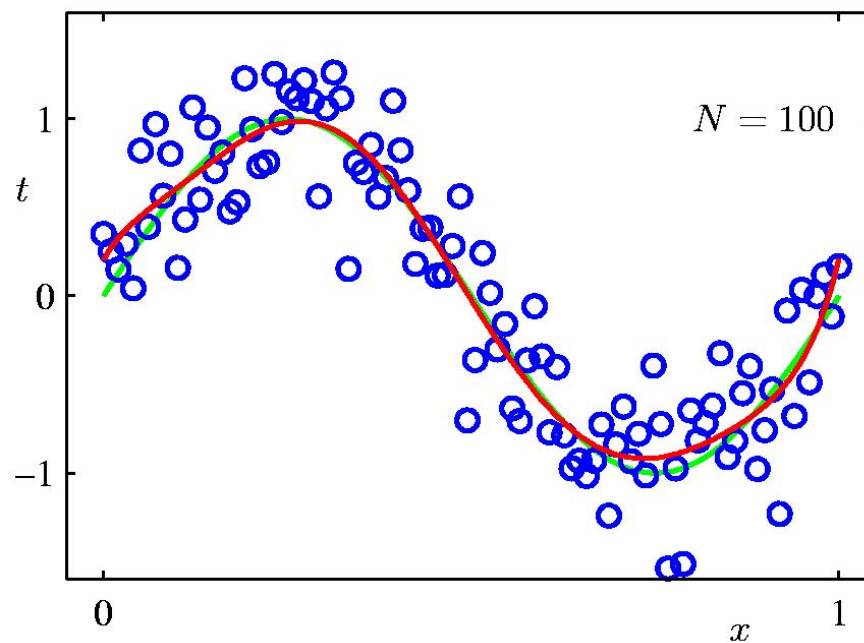
Root-Mean-Square (RMS) Error vs Polynomial order

Data Set Size:

9th Order Polynomial



$N = 15$



$N = 100$

Regularization (to avoid overfitting)

- “regularization term” imposes penalty on less desirable solutions
 - $\text{Cost} = \text{MSE} + \lambda \text{ Penalty}(f)$
 - Regularization Penalty is a functional (maps each function f onto a number)
 - Popular Penalties
 - **ridge regression** (sum squared of weights)
 - **Lasso** (sum of $|w|$; for large λ yields sparse models)
 - **Elastic net**: combines both ridge and Lasso
 - number of non-zero weights
 - smoothness of function
- (**note**: “intercept”, i.e. w_0 , not included in penalties)

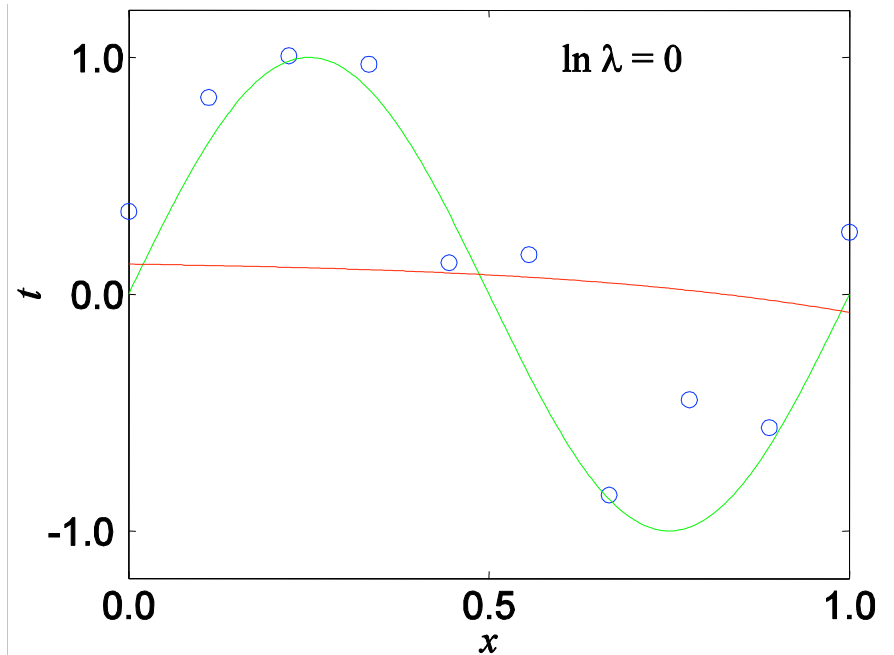
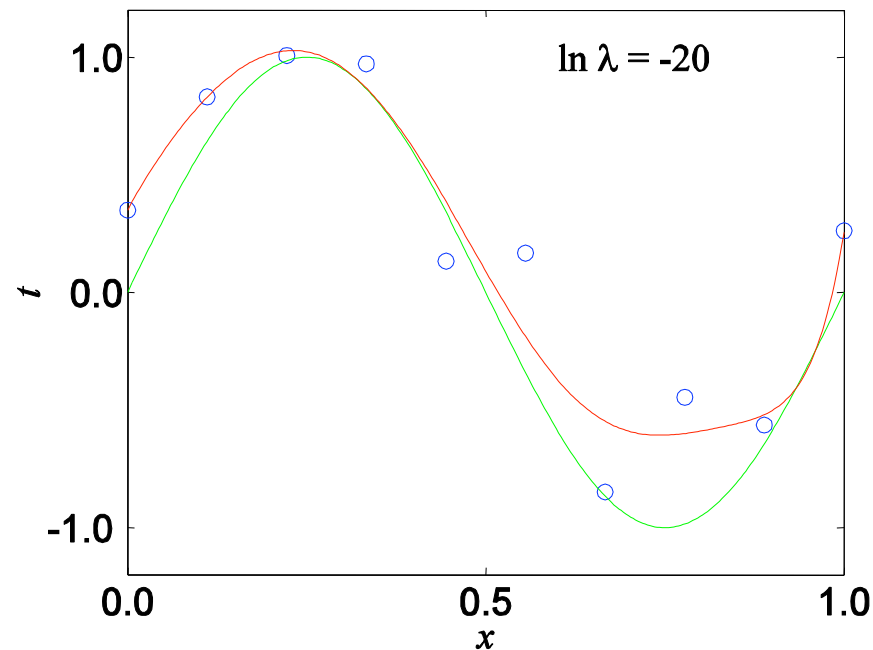
Ridge Regression Example

- Discourage large values by adding penalty term to error

$$E(\mathbf{w}) = \sum_{n=1}^N \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- Also called *shrinkage* (stats) or *weight decay* (neural nets)
- The regularization coefficient λ now controls the effective model complexity
- *Closed form solution: $\mathbf{w} = \left(\lambda \mathbf{I} + \Phi^T \Phi \right)^{-1} \Phi^T \mathbf{t}$.
 - Leads to numerical stability as well!
- *(Related to MAP estimate in Bayesian Linear regression with isotropic covariances for both likelihood and prior, and with $\lambda = \sigma_{\text{likelihood}}^2 / \sigma_{\text{prior}}^2$)

Regularized $M = 9$ Polynomial



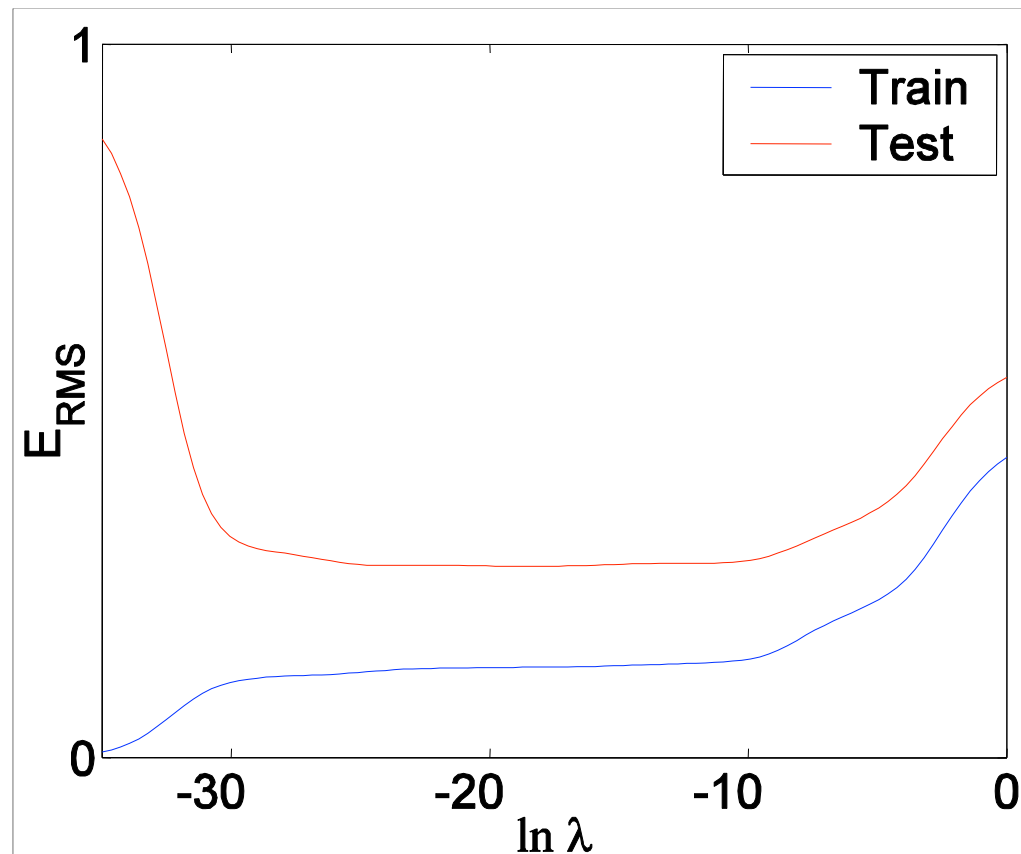
Regularized Parameters

- First col is the unregularized solution

	$\ln \lambda = -\infty$	$\ln \lambda = -20$	$\ln \lambda = 0$
w_0^*	0.35	0.35	0.1273
w_1^*	232.37	5.56	-0.0459
w_2^*	-5321.83	-12.27	-0.0578
w_3^*	48568.31	19.01	-0.0460
w_4^*	-231639.30	-82.58	-0.0321
w_5^*	640042.26	46.49	-0.0201
w_6^*	-1061800.52	141.84	-0.0104
w_7^*	1042400.18	-29.57	-0.0028
w_8^*	-557682.99	-231.55	0.0032
w_9^*	125201.43	142.98	0.0080

Generalization

- Noisy sine problem



Evaluation

- Quality criterion for regression
 - Mean squared error (MSE) or equivalent, e.g. SSE, RMSE
 - true vs. empirical
 - normalized (R^2 value = % of variance explained)
 - Adjusted R^2

Estimating True Performance (Formula Driven)*

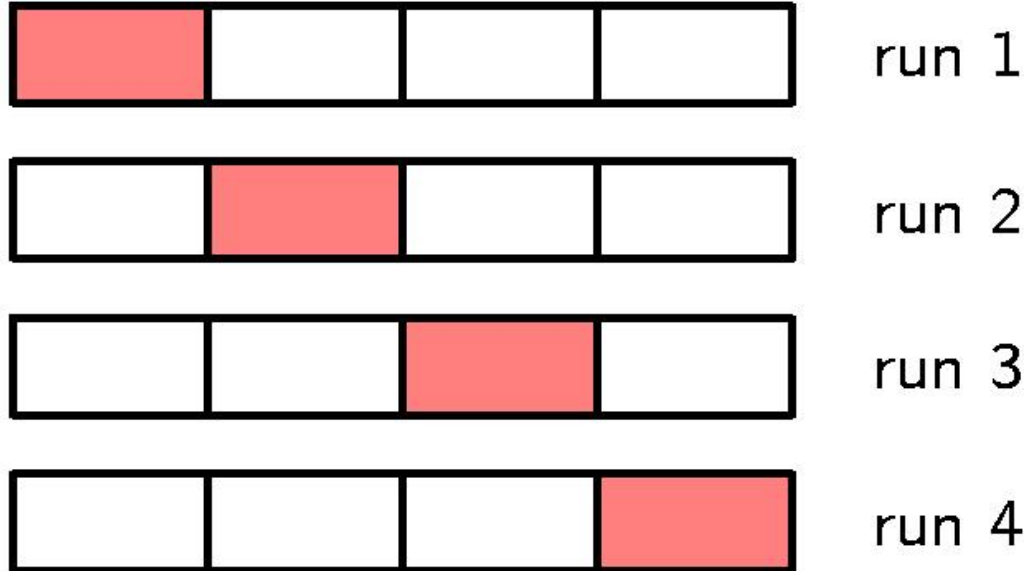
- true mean squared error ($\text{MSE} = \text{SSE}/N$) = empirical error + complexity term
 - complexity term = $f(\text{model type, \# of parameters, \# of training points})$
 - e.g. linear regression with N samples, P parameters
Akaike's Final Prediction error = $\text{MSE}_{\text{empirical}}(N+P) / (N - P)$
 - for nonlinear models, find “effective number of parameters” and plug into linear formulae

Takeaway: Formula Driven Estimates of True Performance specialized for linear models. Not so relevant in data mining context

Estimating True Performance (Data Driven)

- enough data? Use “holdout” to estimate
- Moderately large? Use k-fold cross-validation
 - extreme case (small dataset) : Leave One Out (LOO)

K = 4 example



Bias-Variance Dilemma

Usually *measured* output is not a deterministic function of *given* inputs

Assume: $t = h(\mathbf{x}) + \text{zero-mean noise}$

- your model gives $y(\mathbf{x})$. The *expected squared loss*,

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \iint \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

- **best predictor**: $\mathbb{E}[t \mid \mathbf{x}] = h(\mathbf{x})$;
- $\text{MSE}_{\text{opt}} = \text{variance of the noise inherent in the random variable } t.$
(2nd term on RHS)
- **What does the first term comprise of ?**

The Bias-Variance Decomposition

- Suppose we were given multiple data sets, **each of size N** . Any particular data set, \mathcal{D} , will give a particular function $y(\mathbf{x}; \mathcal{D})$. We then have

$$\begin{aligned} & \{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2 \\ &= \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] + \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ &= \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2 + \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ &\quad + 2\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}. \end{aligned}$$

The Bias-Variance Decomposition (2)

- Taking the expectation over \mathcal{D} yields

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2] \\ = \underbrace{\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2}_{(\text{bias})^2} + \underbrace{\mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2]}_{\text{variance}}. \end{aligned}$$

- (try to express both terms in words)

The Bias-Variance Decomposition (3)

Thus we can write

where $\text{expected loss} = (\text{bias})^2 + \text{variance} + \text{noise}$

$$(\text{bias})^2 = \int \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x}$$

$$\text{variance} = \int \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2] p(\mathbf{x}) d\mathbf{x}$$

$$\text{noise} = \iint \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

Bias: how good the average model is;

Variance: how sensitive the model is to variations in data.

NOTE: the bias and variance concepts here apply to a predictive model, rather than to an estimator of a specific value.

Bias-Variance Tradeoff

- Change model type? Affect bias
- More training data: decrease variance
 - “consistent estimators” converge to ideal solution as $|D| \rightarrow \text{infinity}$
 - For small data sets, lower complexity models may be preferred.
- **Ideal solution: suitable model type & complexity**

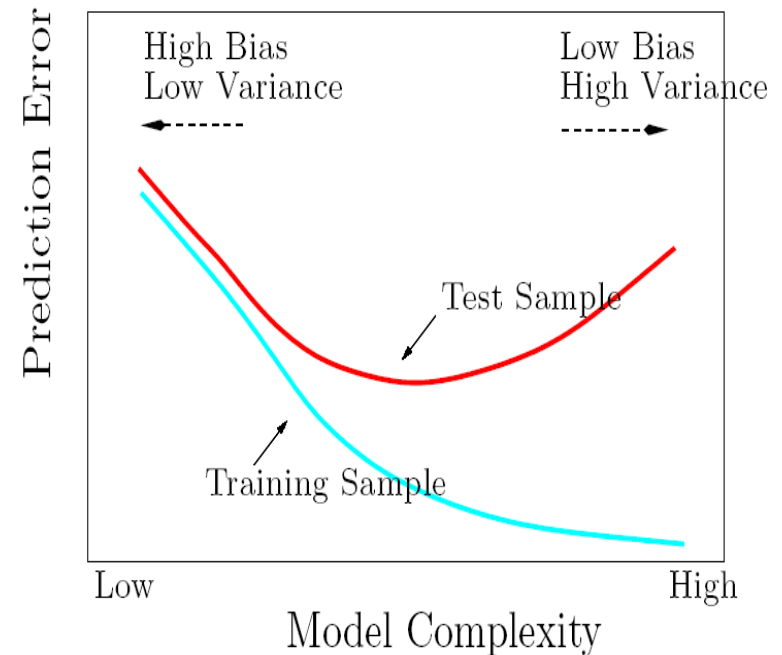
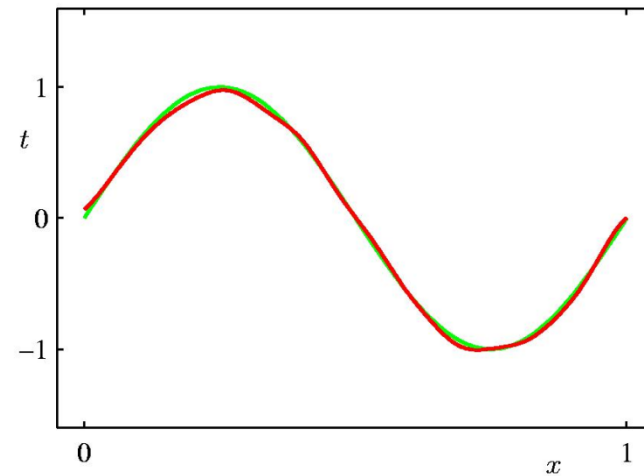
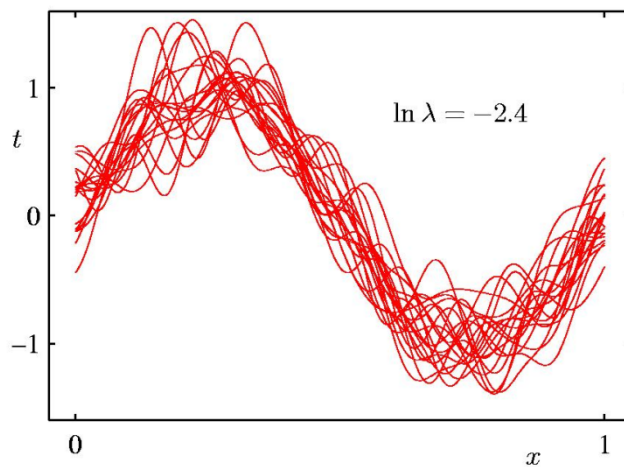
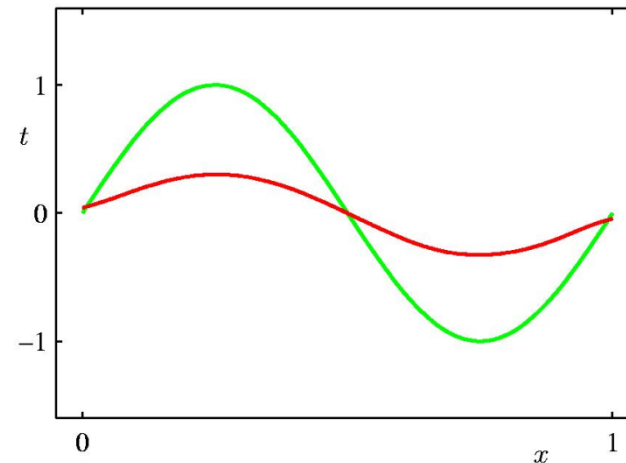
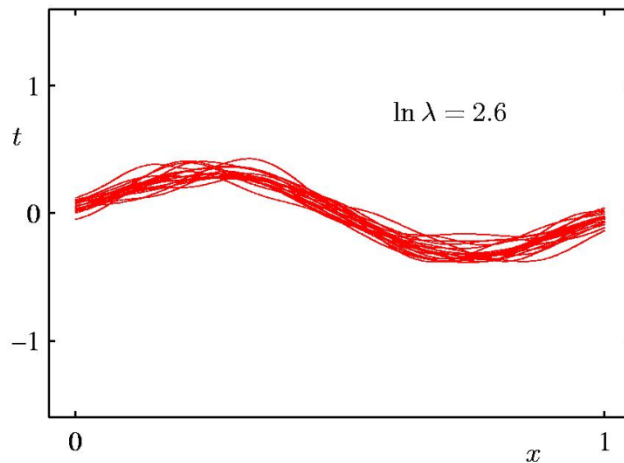


Figure 2.11: *Test and training error as a function of model complexity.*

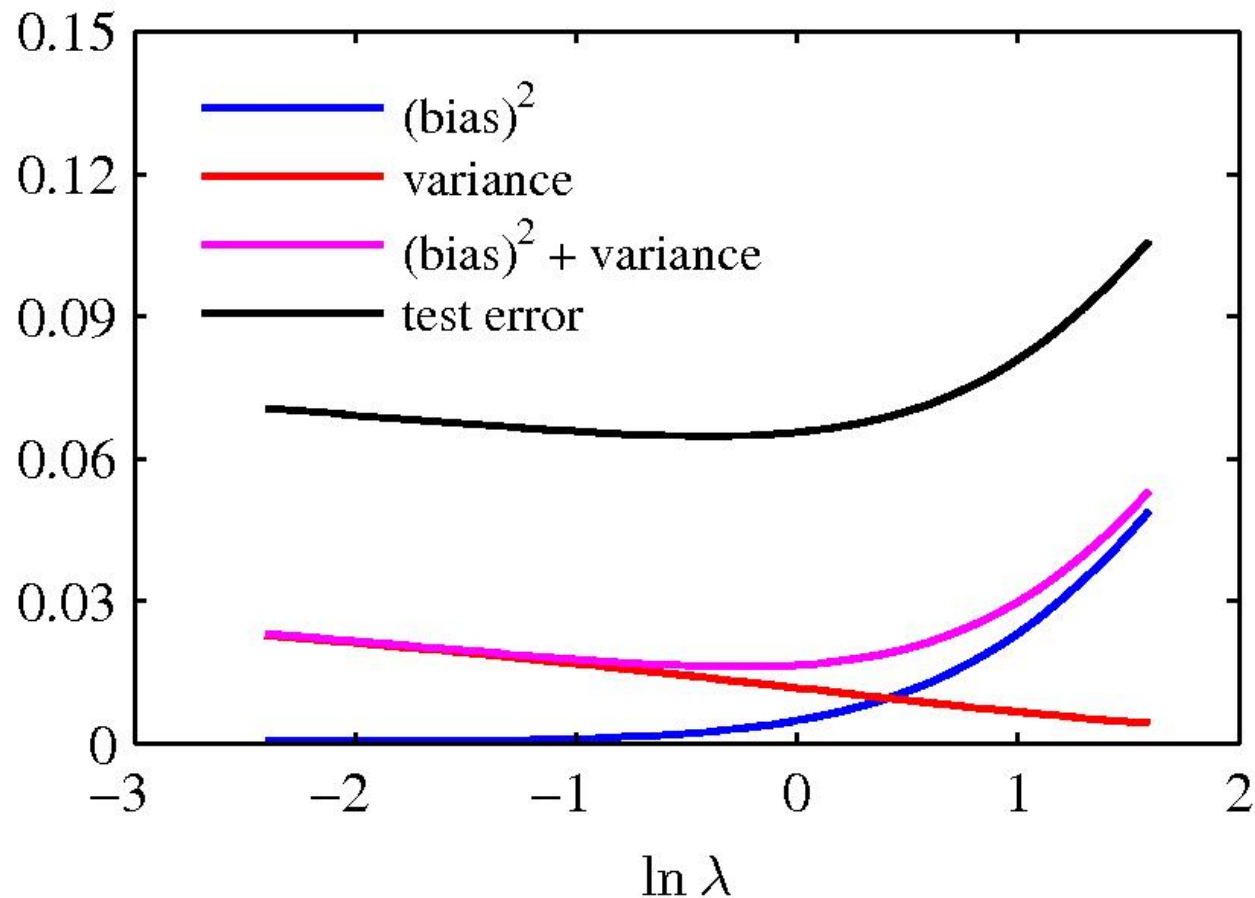
Effect of Regularization on Bias-Variance

- Bishop 06, fig 3.5. Model is sum of 24 gaussians, with ridge regression



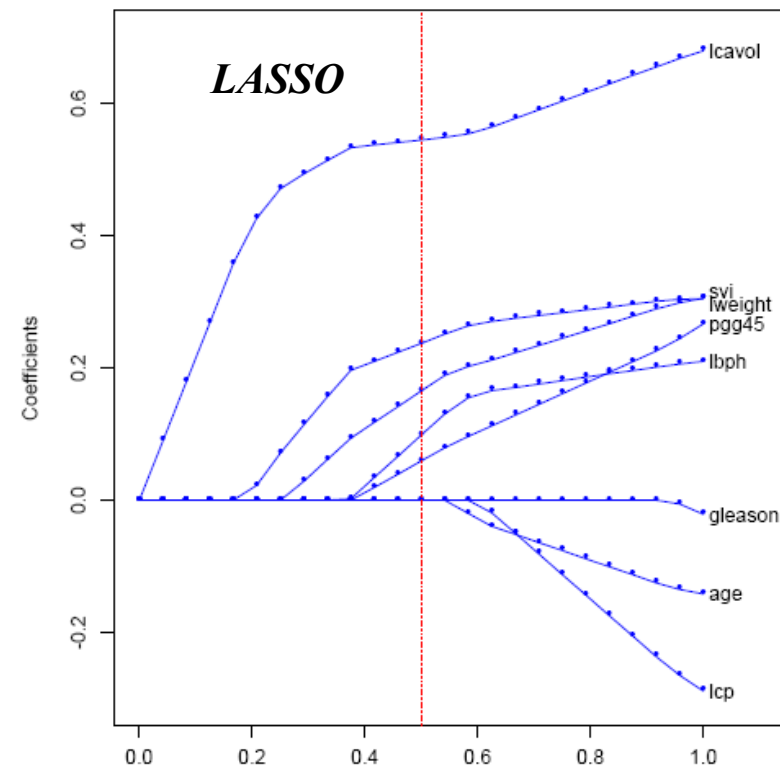
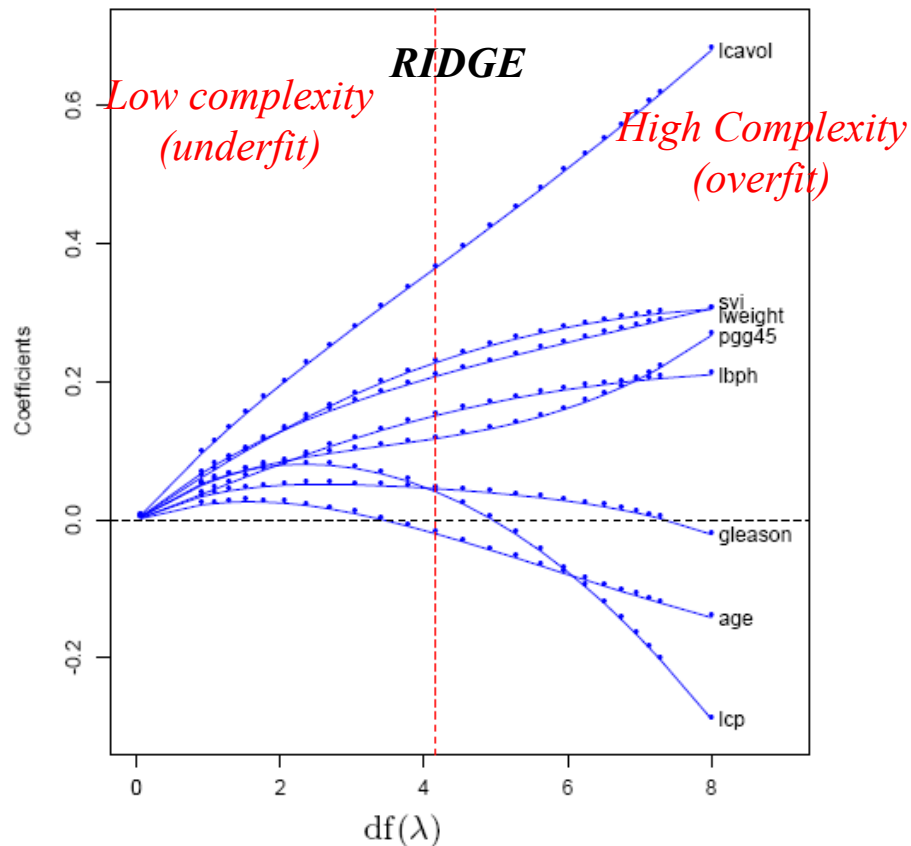
Bias-Variance vs. Regularization Amount

- What happens to the curves as amount of training data increases?*



Ridge vs. Lasso

- HTF figs 3.7, 3.9: Prostate Cancer example. Red line chosen by Cross-validation



Effect on values of coefficients as “effective degrees of freedom (DOF)” is increased for

(a) Ridge regression (left) and (b) Lasso (Right).

High λ translates to low DOF, so λ is being progressively decreased from left to right along the x-axis.


Extras

Geometry of Least Squares*

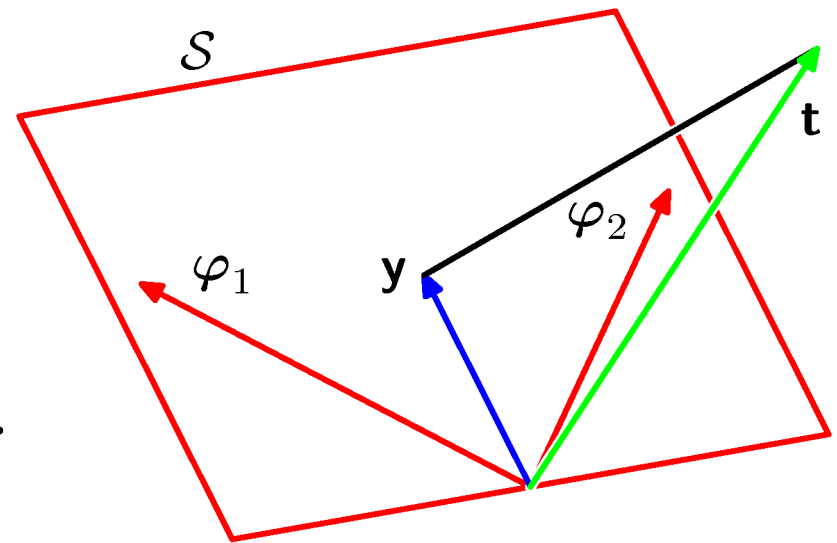
- Consider

$$\mathbf{y} = \Phi \mathbf{w}_{\text{ML}} = [\varphi_1, \dots, \varphi_M] \mathbf{w}_{\text{ML}}.$$

$$\mathbf{y} \in \mathcal{S} \subseteq \mathcal{T} \quad \mathbf{t} \in \mathcal{T}$$


N-dimensional
M-dimensional

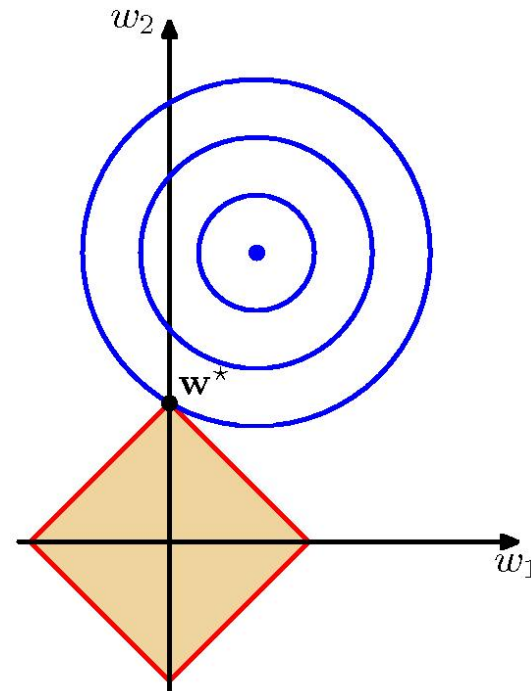
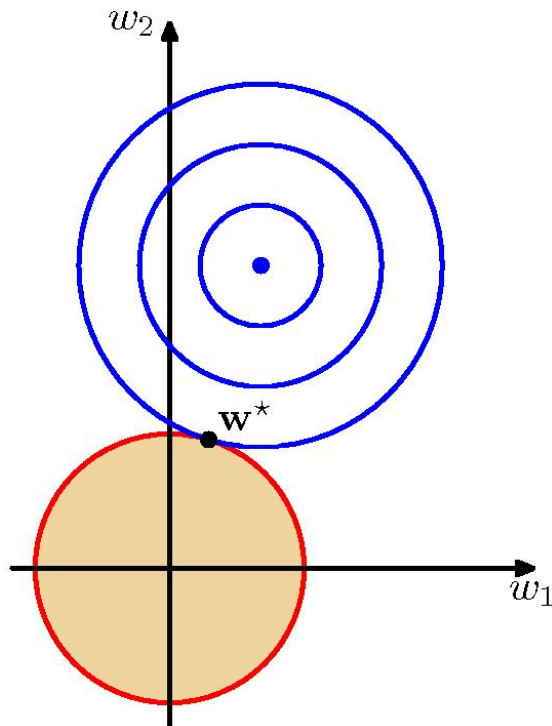
- \mathcal{S} is spanned by $\varphi_1, \dots, \varphi_M$.
- \mathbf{w}_{ML} minimizes the distance between \mathbf{t} and its orthogonal projection on \mathcal{S} , i.e. \mathbf{y} .



Takeaway: You are restricted by your choice of the features

Comparing Shrinkage Methods B06: fig 3.4

- ridge regression (Regularization Penalty = sum squared of weights)
vs
- **Lasso** ((Regularization Penalty = sum of $|w|$)
red: constant penalty contour; blue: unregularized error contours



Least Angle Regression (LAR; HTF 3.4.4)*

- Takeaway: Efficient procedure for fitting an entire lasso sequence with the cost of a single least squares fit.

- R code: `lar`

Algorithm 3.2 *Least Angle Regression.*

1. Standardize the predictors to have mean zero and unit norm. Start with the residual $\mathbf{r} = \mathbf{y} - \hat{\mathbf{y}}$, $\beta_1, \beta_2, \dots, \beta_p = 0$.
 2. Find the predictor \mathbf{x}_j most correlated with \mathbf{r} .
 3. Move β_j from 0 towards its least-squares coefficient $\langle \mathbf{x}_j, \mathbf{r} \rangle$, until some other competitor \mathbf{x}_k has as much correlation with the current residual as does \mathbf{x}_j .
 4. Move β_j and β_k in the direction defined by their joint least squares coefficient of the current residual on $(\mathbf{x}_j, \mathbf{x}_k)$, until some other competitor \mathbf{x}_l has as much correlation with the current residual.
 5. Continue in this way until all p predictors have been entered. After $\min(N - 1, p)$ steps, we arrive at the full least-squares solution.
-

Group Lasso for Sparse Learning*

- SLEP package
<http://www.public.asu.edu/~jye02/Software/SLEP/overview.htm>
- ℓ_1 -Regularized (Constrained) Sparse Learning
- ℓ_1/ℓ_q -Regularized Sparse Learning ($q>1$)
- Fused Lasso
- Sparse Inverse Covariance Estimation
- Sparse Group Lasso
- Tree Structured Group Lasso
- Overlapping Group Lasso
- *Takeaway: A variety of methods exist to shrink parameters in different ways*