
Function approximation (Regression)

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".

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)

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

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

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

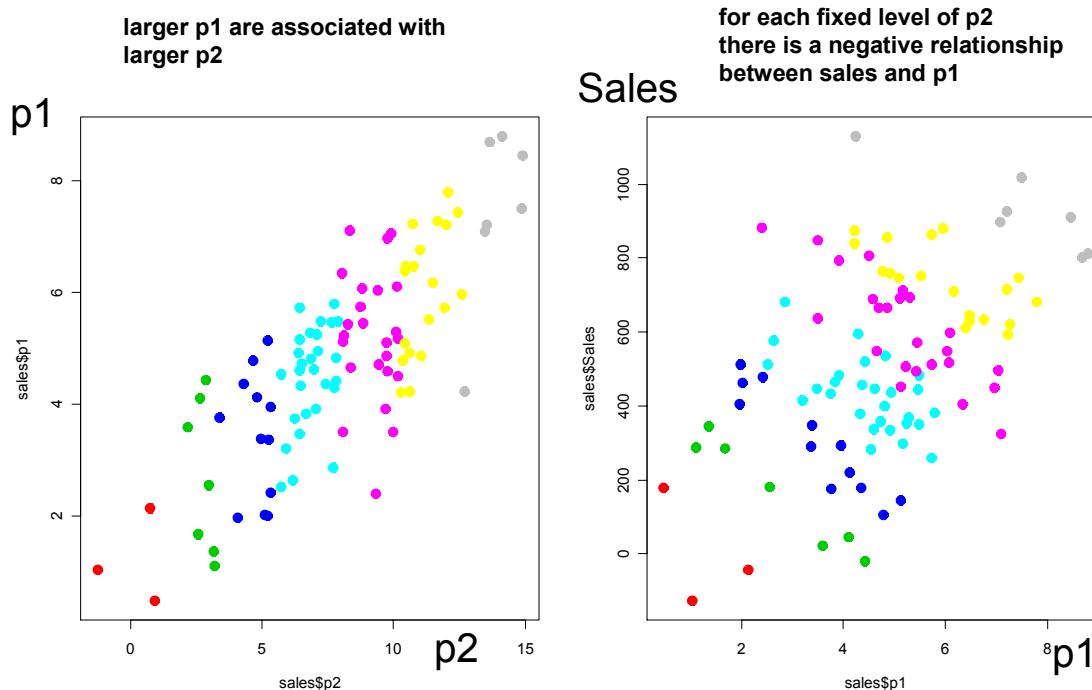
so that the basis functions are given by $\phi_i(x) = x^i$

Assumptions

- Expected value of T given the basis function vector phi is linear in phi.
 - Conditional mean is linear in the predictors.
- All distributions around the expected values are assumed to be i.i.d. zero mean Gaussian with constant variance.
 - In stats notation: $Y|X_1 \dots X_p \sim N(\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p, \sigma^2)$
 - In new notation: (fill in)
- How can you “verify” that the assumptions seem reasonable?
- Consequences of the assumptions?
 - minimizing Mean Squared Error (MSE) on the training data yields the Maximum Likelihood Estimate (MLE) solution of the assumed **generative** model. (Bishop pg. 140)

How do you interpret the weights?

- What does a partial derivative imply?
- Example, my sales (t) as a function of my price (p_1) and competitor's price (p_2)
 - See dataset on Canvas.

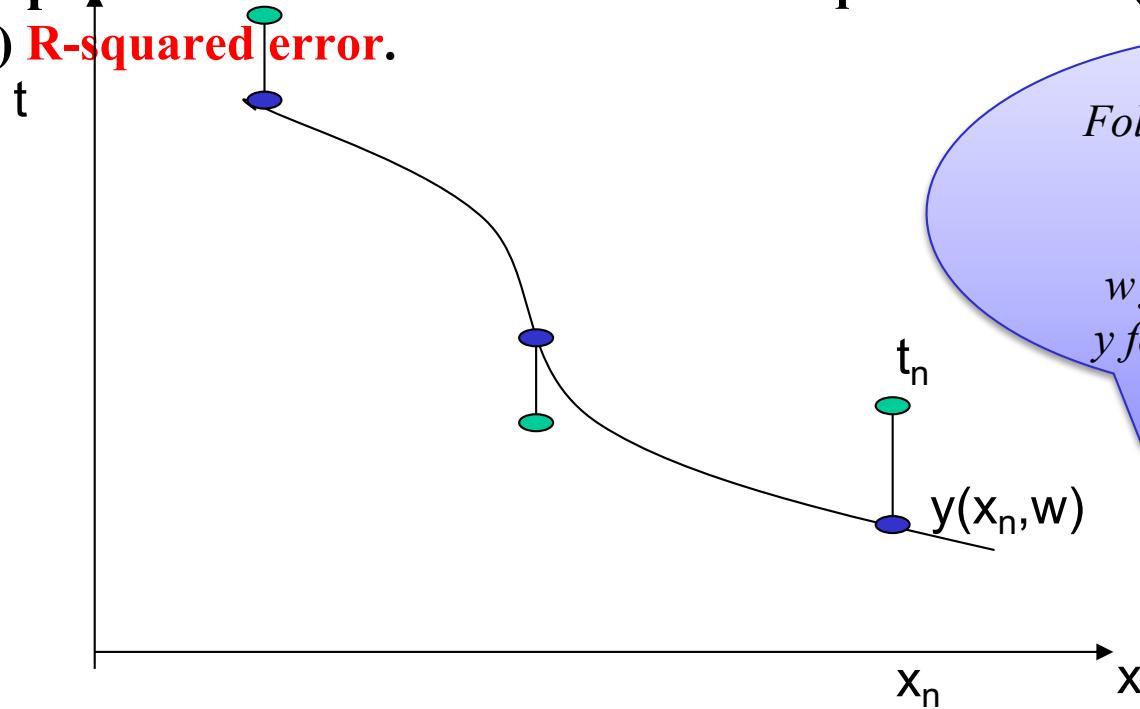


Ordinary Least Squares (OLS)

- Minimize a **loss function** $E(w)$ given by sum-of-squares **error (SSE)**
(t 's are the target values)

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

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)

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

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

- “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 an $(M+1) \times (M+1)$ matrix

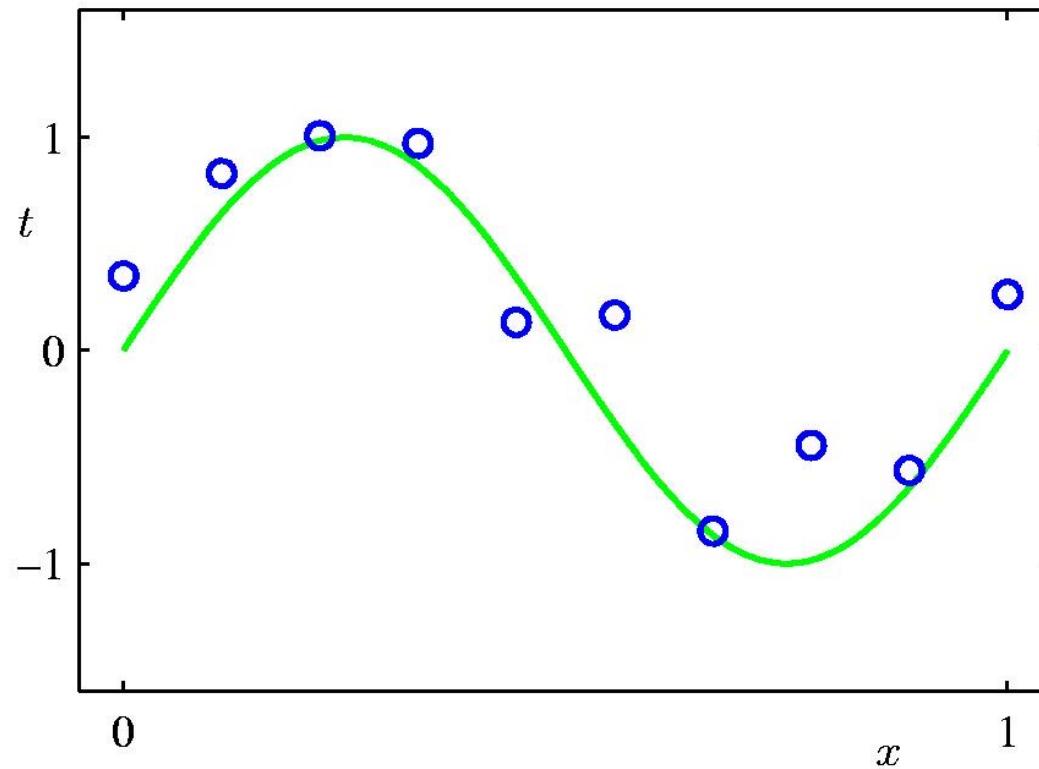
Explicitly shows collinearity problem

Collinearity problem: parameter estimates have high uncertainty if two or more independent variables are highly collinear

in MLR, the standard errors are defined by the following formula:

$$s_{b_j}^2 = \frac{s^2}{(N-1)(\text{variation in } X_j \text{ not associated with other } X's)}$$

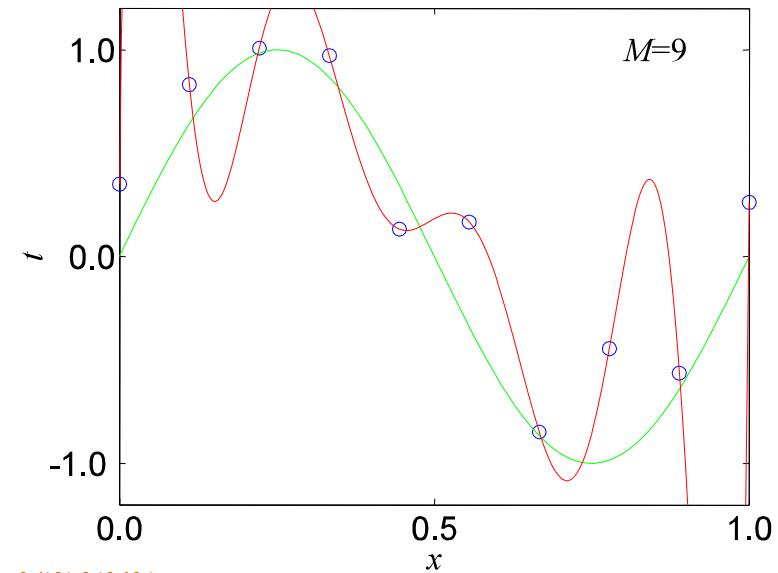
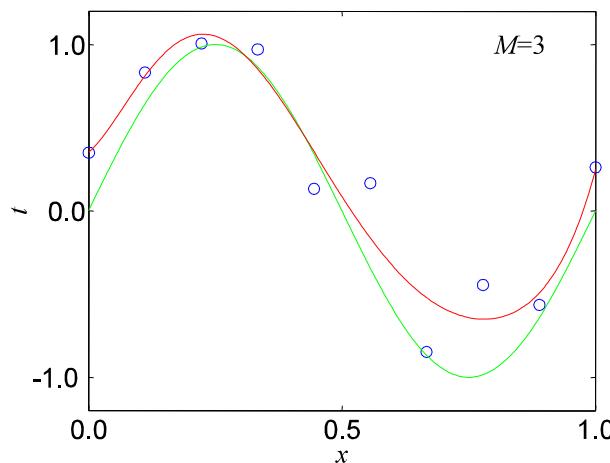
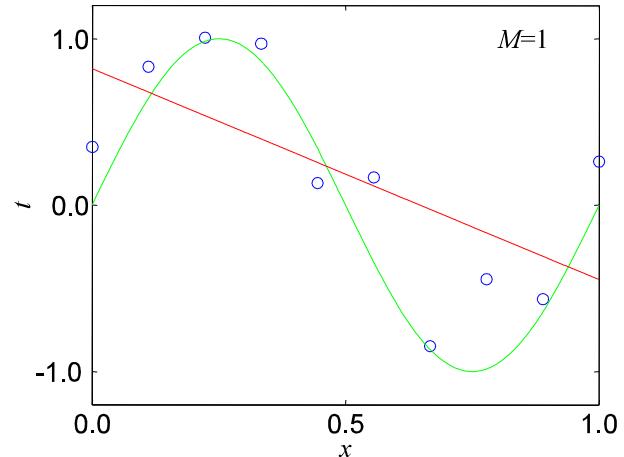
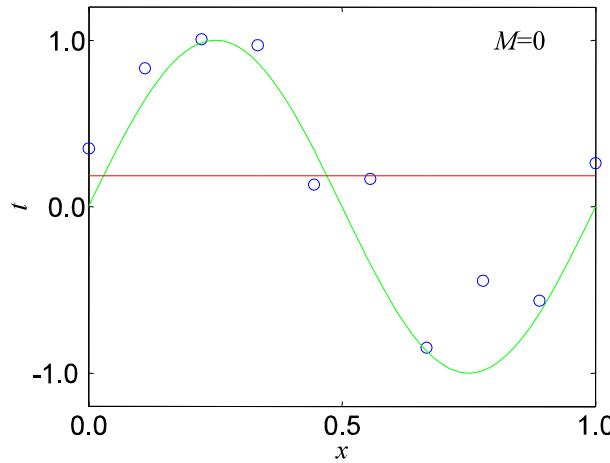
Polynomial Curve Fitting



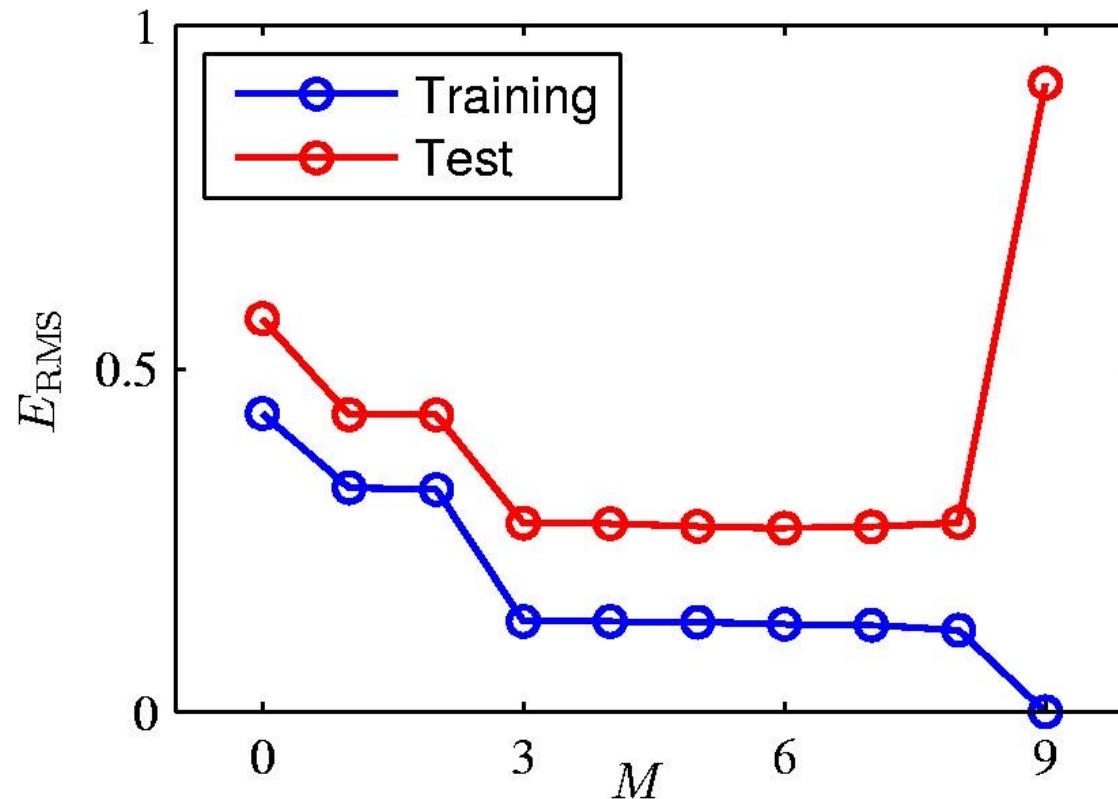
$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$

Model Complexity and Overfitting

- “Noisy sine” example from Chris Bishop



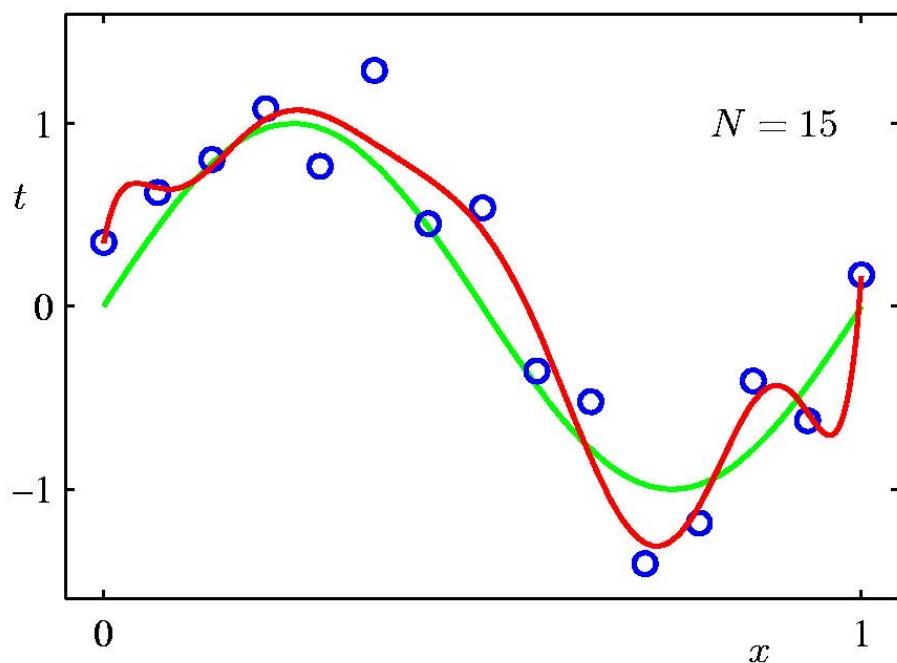
Over-fitting



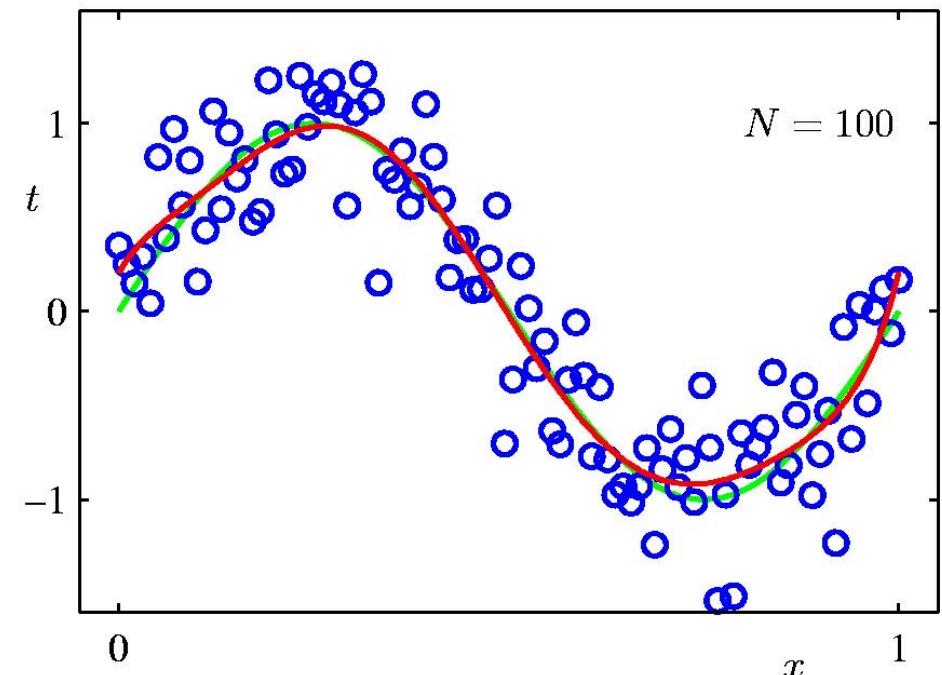
Root-Mean-Square (RMS) Error vs Polynomial order

Data Set Size:

9th Order Polynomial



$N = 15$



$N = 100$

Regularization (to avoid overfitting)

- “regularization term” imposes penalty on less desirable solutions
 - Cost = MSE + λ Penalty (f)
 - Scikit uses “alpha” to denote lambda.
 - 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:** 1. “intercept”, i.e. w_0 , not included in penalties
2. customary to standardize all independent variables first (why?)

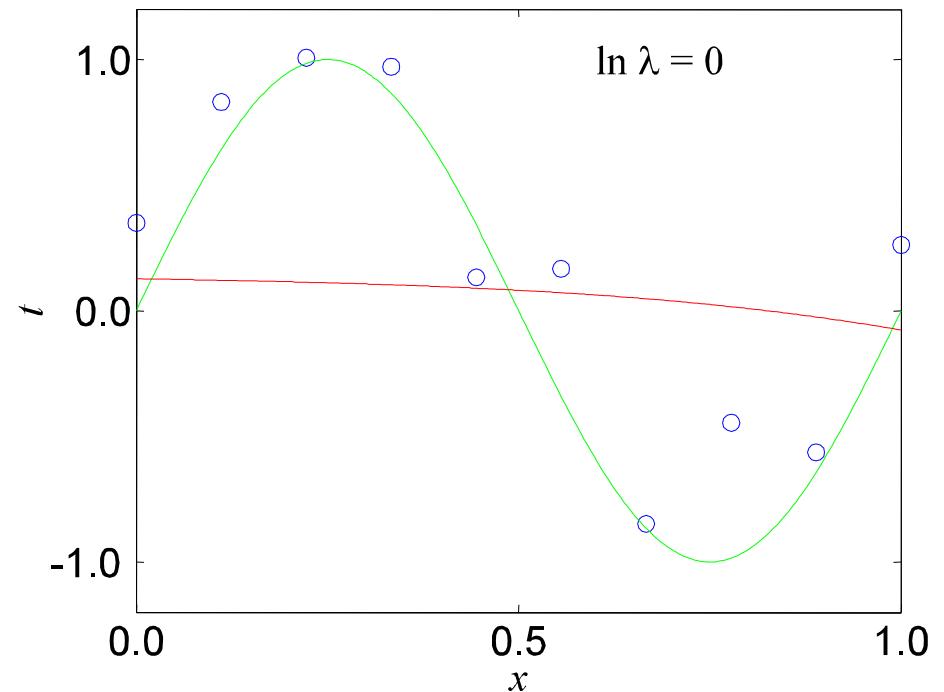
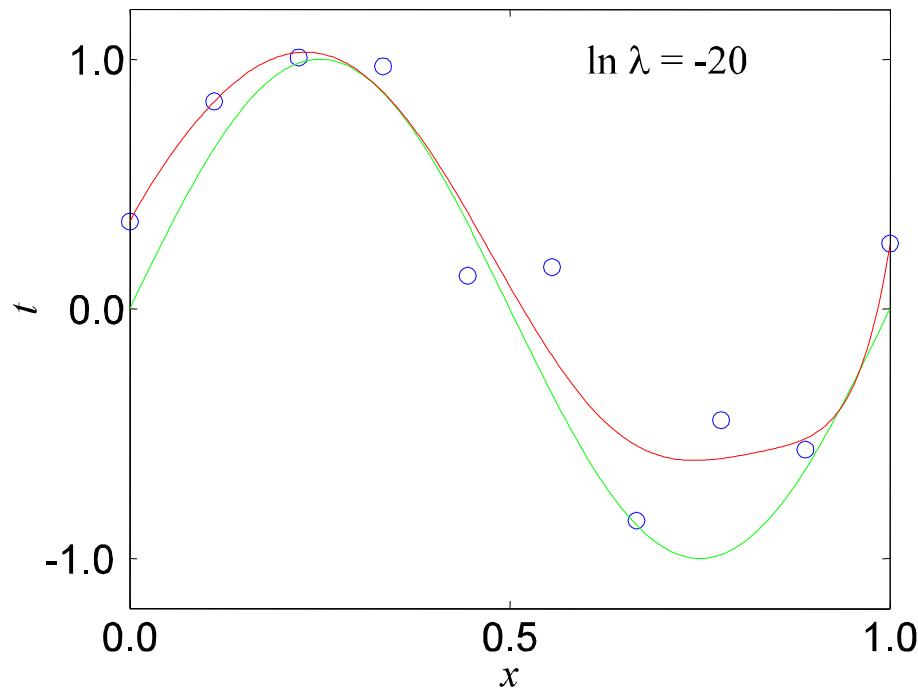
Ridge Regression Example

- Discourage large values by adding penalty term to error

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

- Also called *shrinkage* (stats) or *weight decay (neural nets)*
- The regularization coefficient λ now controls the effective model complexity
- *Closed form solution: $\mathbf{w} = (\lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{t}$.
 - Leads to numerical stability as well!

Regularized $M=9$ Polynomial



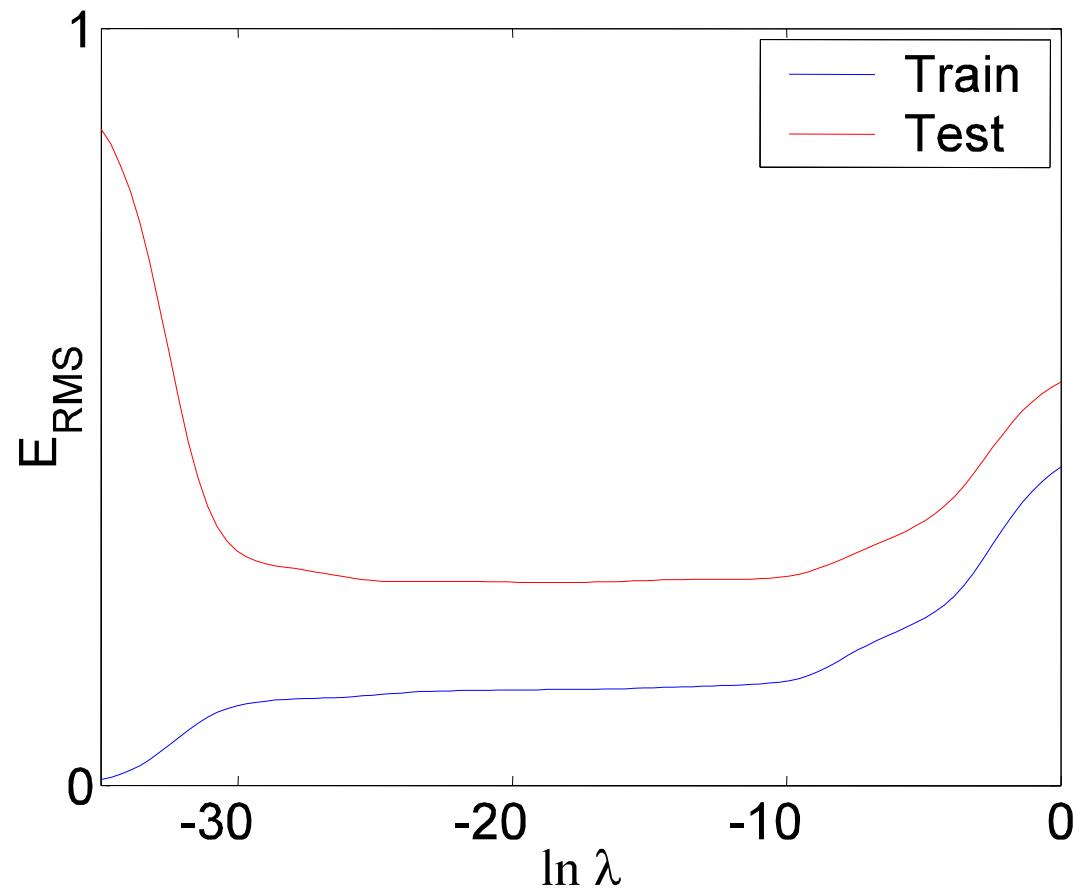
Regularized Parameters

- 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

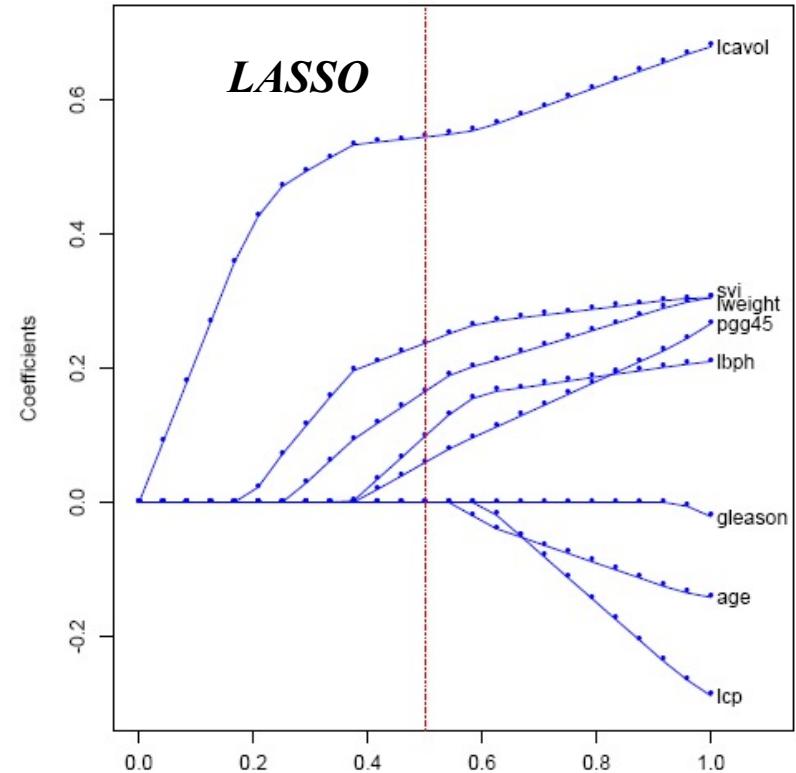
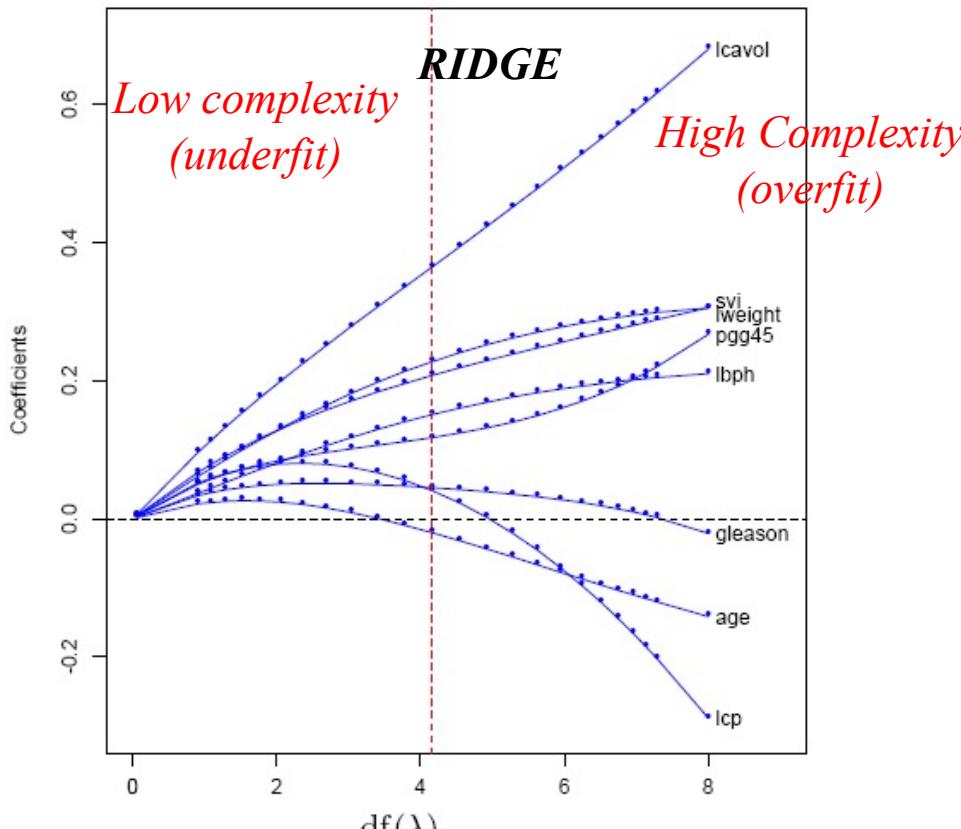
Generalization

- Noisy sine problem



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 (df)” is increased for
 (a) Ridge regression (left) and (b) Lasso (Right).*

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

Evaluation

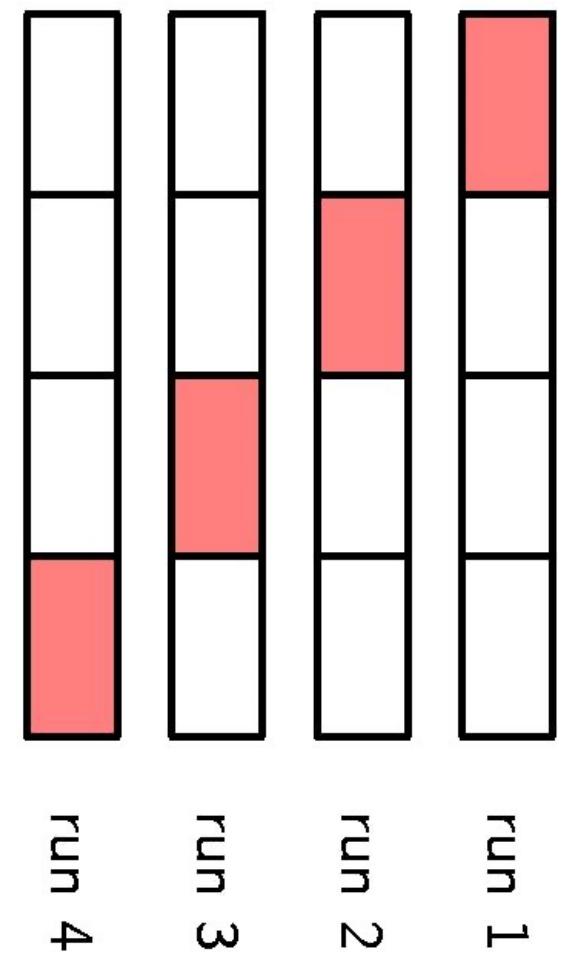
- Quality criterion for regression
 - Mean squared error (MSE) or equivalent, e.g. SSE, RMSE
 - true vs. empirical
 - normalized (R^2 value = % of variance explained)
 - Adjusted R^2

Estimating True Performance (Data Driven)

- enough data? Use “holdout” to estimate
- Moderately large? Use k-fold cross-validation

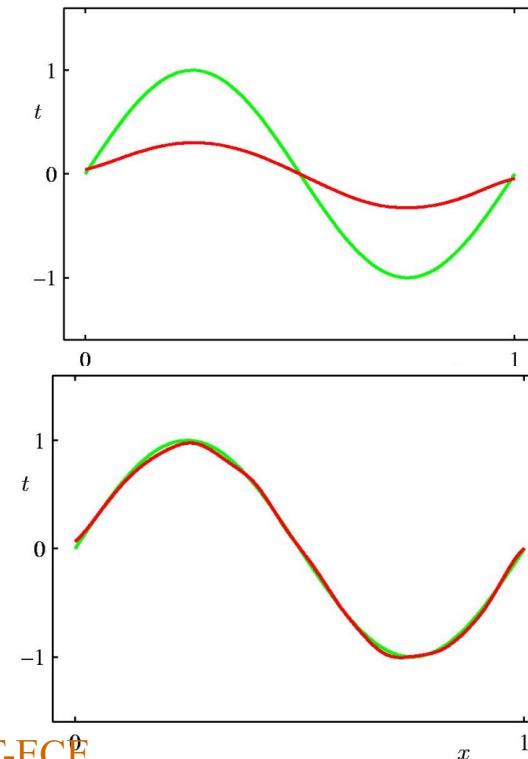
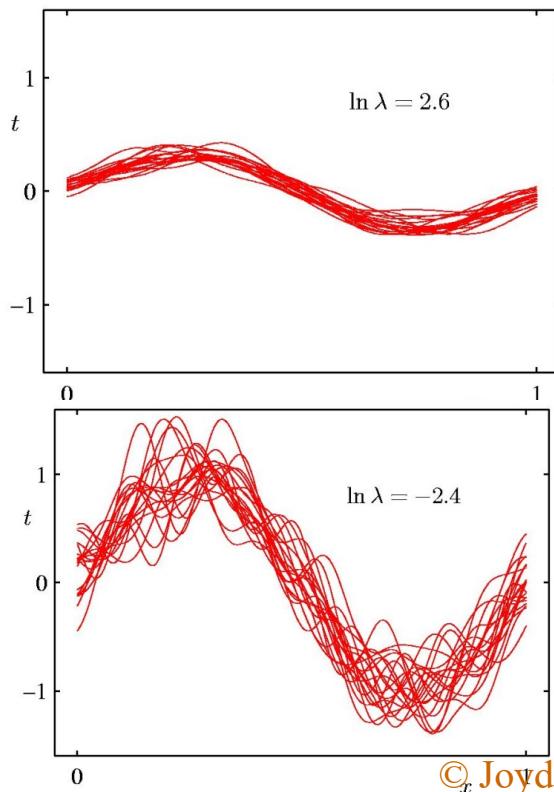
$k = 4$ example

Shaded subset: validation set for that run



How to evaluate a Regression Model?

- Model = math form + learning method.
 - Will be impacted by given training (and validation) dataset!!
 - Want to evaluate the model irrespective of the specific train/validate dataset used.
 - Need to consider the collection of solutions obtained, not one specific solution.
- Bishop 06, fig 3.5. Compares highly regularized solution (top row) vs. less regularized solution (bottom row). Individual solutions (left), averaged solution (right)



Bias-Variance Dilemma

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

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

- your model gives $y(x)$. 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:** $\mathbb{E}[t | x] = h(x)$;
 - $\text{MSE}_{\text{opt}} = \text{variance of the noise inherent in the random variable } t.$
(2nd term on RHS)
- **What does the first term comprise of ?** (Model_bias)² + Model_variance
 - Bias: how good the average model is;
 - Variance: how sensitive the model is to variations in data.
- Tradeoff between the two terms as function of model complexity
- expected loss = (bias)² + variance + noise

Math Details: 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 $y(\mathbf{x}; \mathcal{D})$.
- For any \mathbf{x} , The expected loss (over datasets of size N) is

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

(try to express both terms in words)

The Bias-Variance Decomposition II*

Considering all possible values of \mathbf{x} , we can write

where expected loss = (bias)² + variance + noise

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

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

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

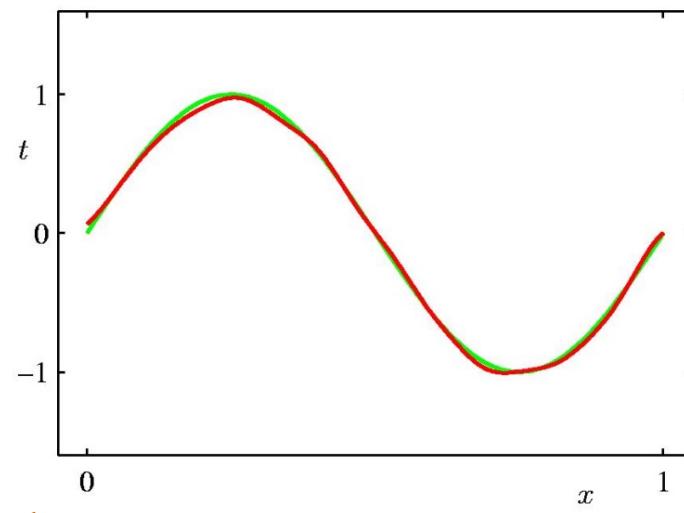
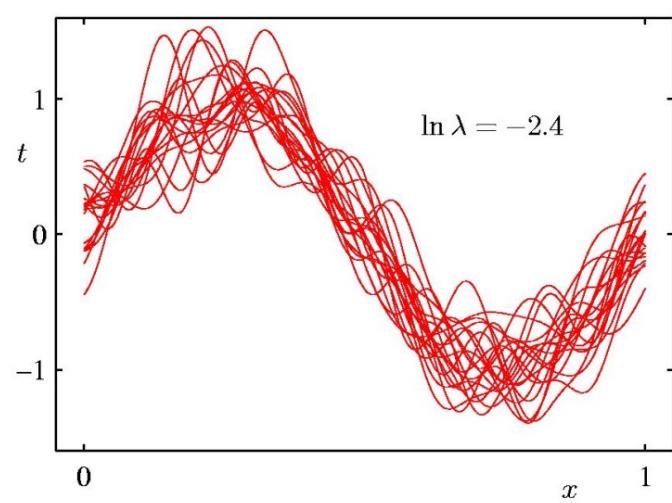
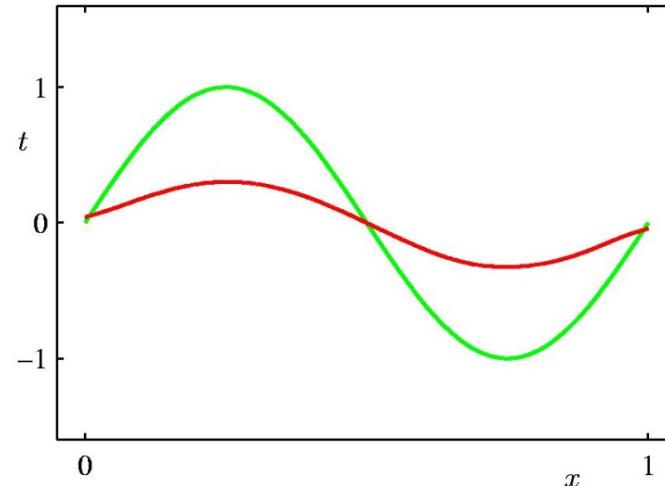
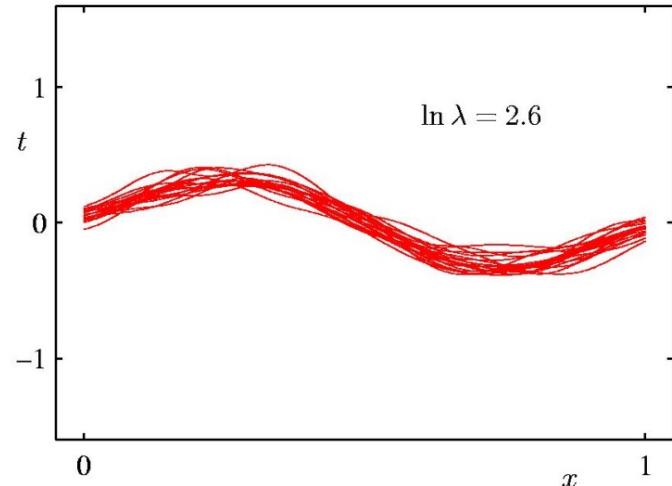
Bias: how good the average model is;

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.

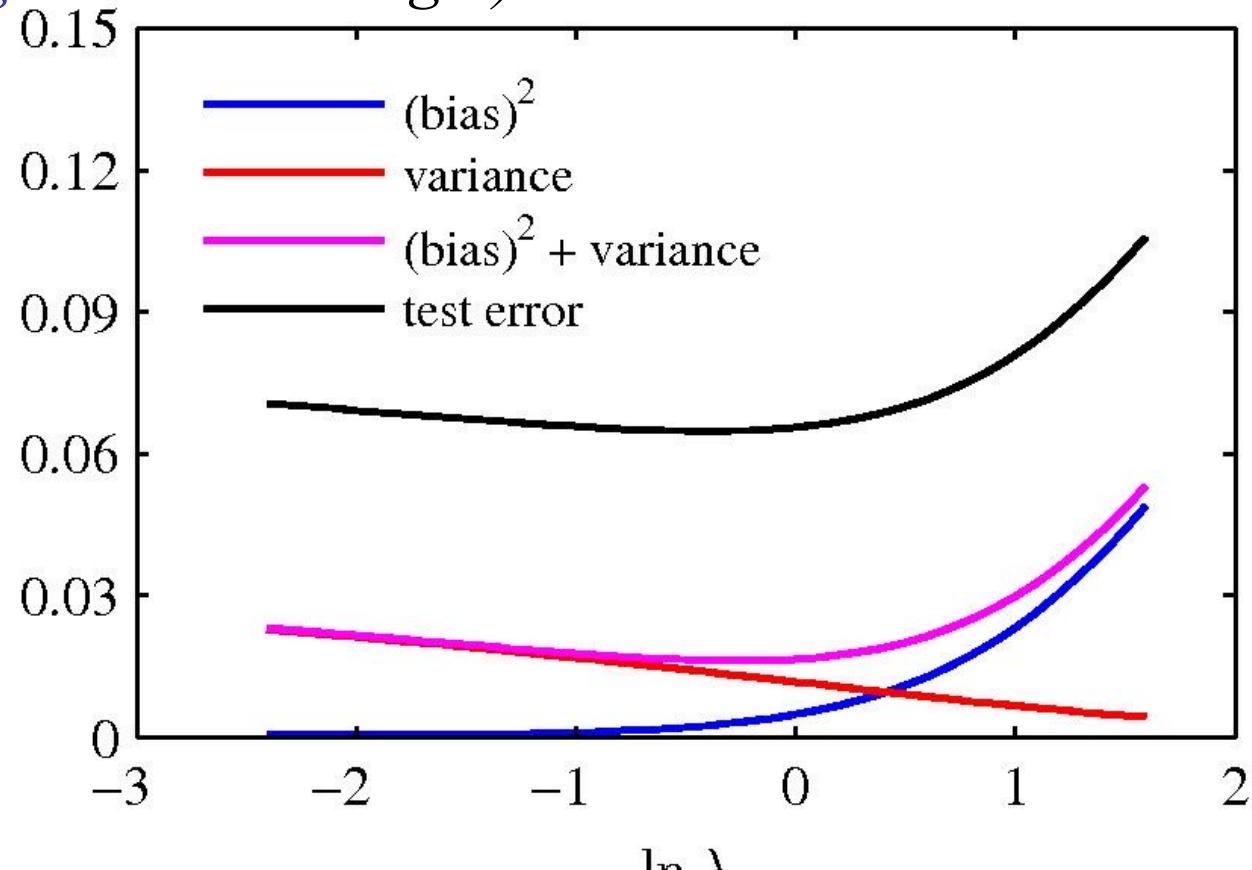
Effect of Regularization on Bias-Variance

- 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 ? (note: effective model complexity is decreasing towards the right)*



Bias-Variance Tradeoff

- **Your task:** qualitatively plot bias² and variance in fig 2.11
- Change model type? Affect bias
- More training data: decrease variance
 - “consistent estimators” converge to ideal solution as $|D| \rightarrow \infty$
 - For small data sets, lower complexity models may be preferred.
- **Ideal solution:** suitable model type & complexity

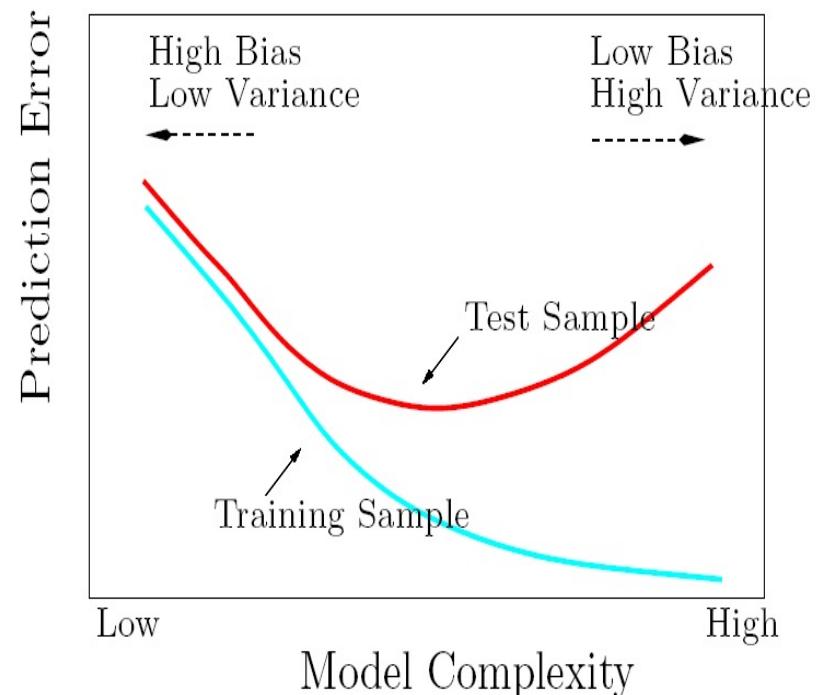
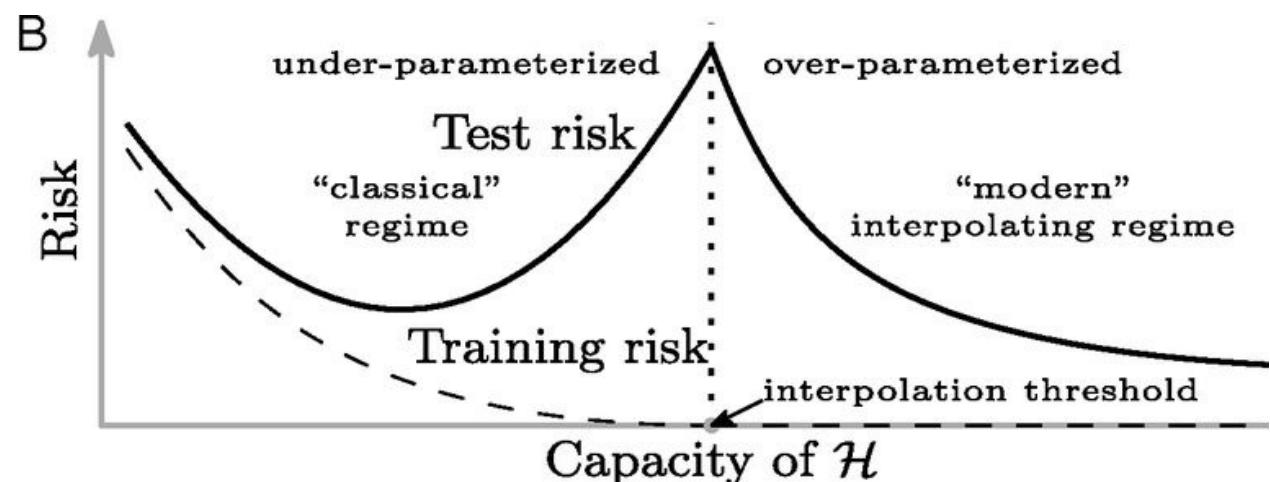
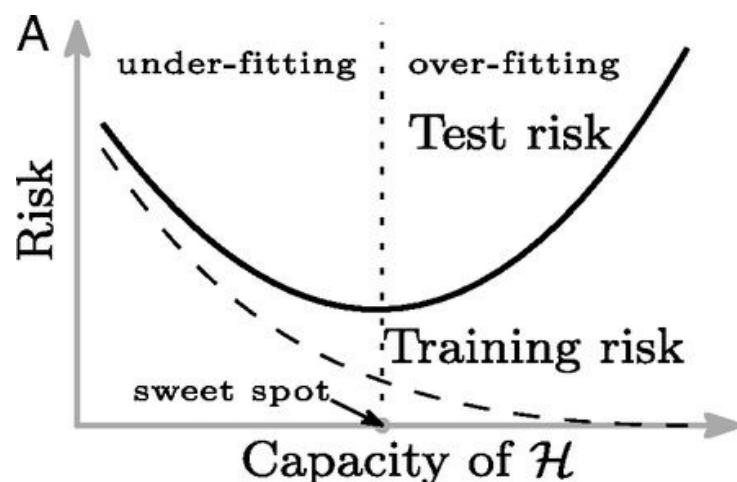


Figure 2.11: *Test and training error as a function of model complexity.*

Breaking News (2019)

- Variance can actually go down in the highly over-parameterized regime (going beyond interpolation)!!
 - Figure below from M. Belkin, D. Hsu, S. Ma, and S. Mandal, “Reconciling modern machine learning practice and the bias-variance trade-off,” *Proceedings of the National Academy of Sciences*, vol. 116, no. 32, pp. 15849–15854, 2019.



Application: How do you improve your model?

- Get more training data
 - Change complexity (e.g. via regularization)
 - Change optimization method
 - Change Model type
-
- Still not acceptable?
 - Change feature space

Bias –variance tradeoff is encountered in many situations

Example: determining # of bins for a histogram.

Extras

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 T
 - Think of training/test datasets as i.i.d. samples from an underlying joint distribution $p(X, T)$
 - Build: a model relating X to T
 - single value for T given X (most common)
 - e.g. $E[T | X]$, the “regression of t on X .
 - Assumes $T = \text{function of } X + (\text{zero-mean, symmetric}) \text{ noise}$
 - Add Confidence Interval (e.g. based on the Normally distributed noise term in MLR)
 - (Arbitrary) Distribution of T given X

Geometry of Least Squares*

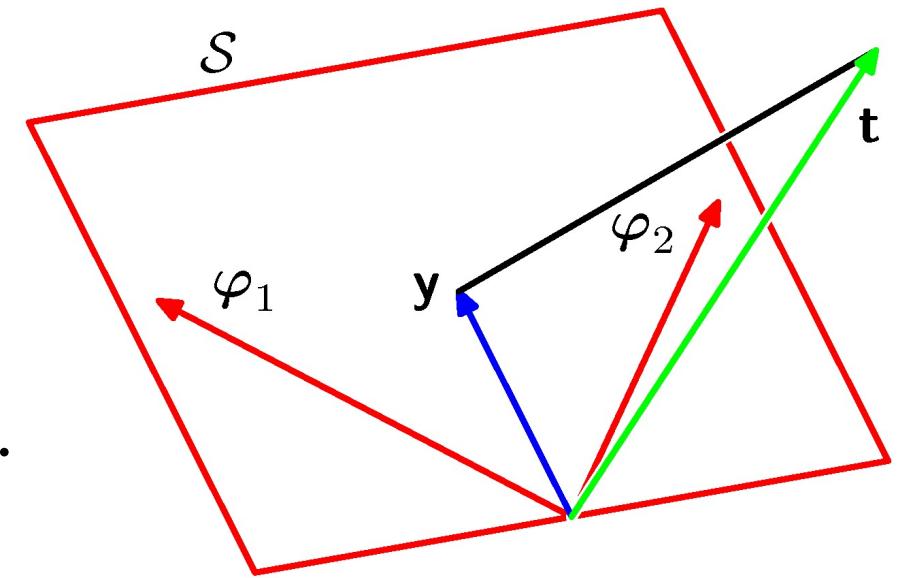
- Consider

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

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

\mathcal{S} : *N-dimensional*
 \mathcal{T} : *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

Correlation vs. Causation

- Though it may still be actionable

Top 10 Best (and Worst) Educated States, and How They Voted

ranked by percentage of residents 25 years of age or older with college degree or more

% over 25 with
college degree

Best Educated

39.1%	1. Massachusetts
36.9%	2. Maryland
36.7%	3. Colorado
36.2%	4. Connecticut
35.4%	5. Vermont
35.3%	6. New Jersey
35.1%	7. Virginia
33.4%	8. New Hampshire
32.9%	9. New York
32.4%	10. Minnesota



% over 25 with
college degree

Worst Educated

18.5%	1. West Virginia
19.8%	2. Mississippi
20.3%	3. Arkansas
21.1%	4. Kentucky
21.1%	5. Louisiana
22.3%	6. Alabama
22.5%	7. Nevada
23.0%	8. Indiana
23.6%	9. Tennessee
23.8%	10. Oklahoma



Research Statistics provided by FoxBusiness.com, based on education data from the U.S. Census Bureau's American Community Survey. 24/7 Wall St. identified the U.S. states with the largest and smallest percentages of residents 25 or older with a college degree or more. <http://www.foxbusiness.com/personal-finance/2012/10/15/americas-best-and-worst-educated-states/>

Incremental Forward Stagewise Regression (HTF 3.8)

Algorithm 3.4 *Incremental Forward Stagewise Regression— FS_ϵ .*

1. Start with the residual \mathbf{r} equal to \mathbf{y} and $\beta_1, \beta_2, \dots, \beta_p = 0$. All the predictors are standardized to have mean zero and unit norm.
 2. Find the predictor \mathbf{x}_j most correlated with \mathbf{r}
 3. Update $\beta_j \leftarrow \beta_j + \delta_j$, where $\delta_j = \epsilon \cdot \text{sign}[\langle \mathbf{x}_j, \mathbf{r} \rangle]$ and $\epsilon > 0$ is a small step size, and set $\mathbf{r} \leftarrow \mathbf{r} - \delta_j \mathbf{x}_j$.
 4. Repeat steps 2 and 3 many times, until the residuals are uncorrelated with all the predictors.
-

HTF, pg 87

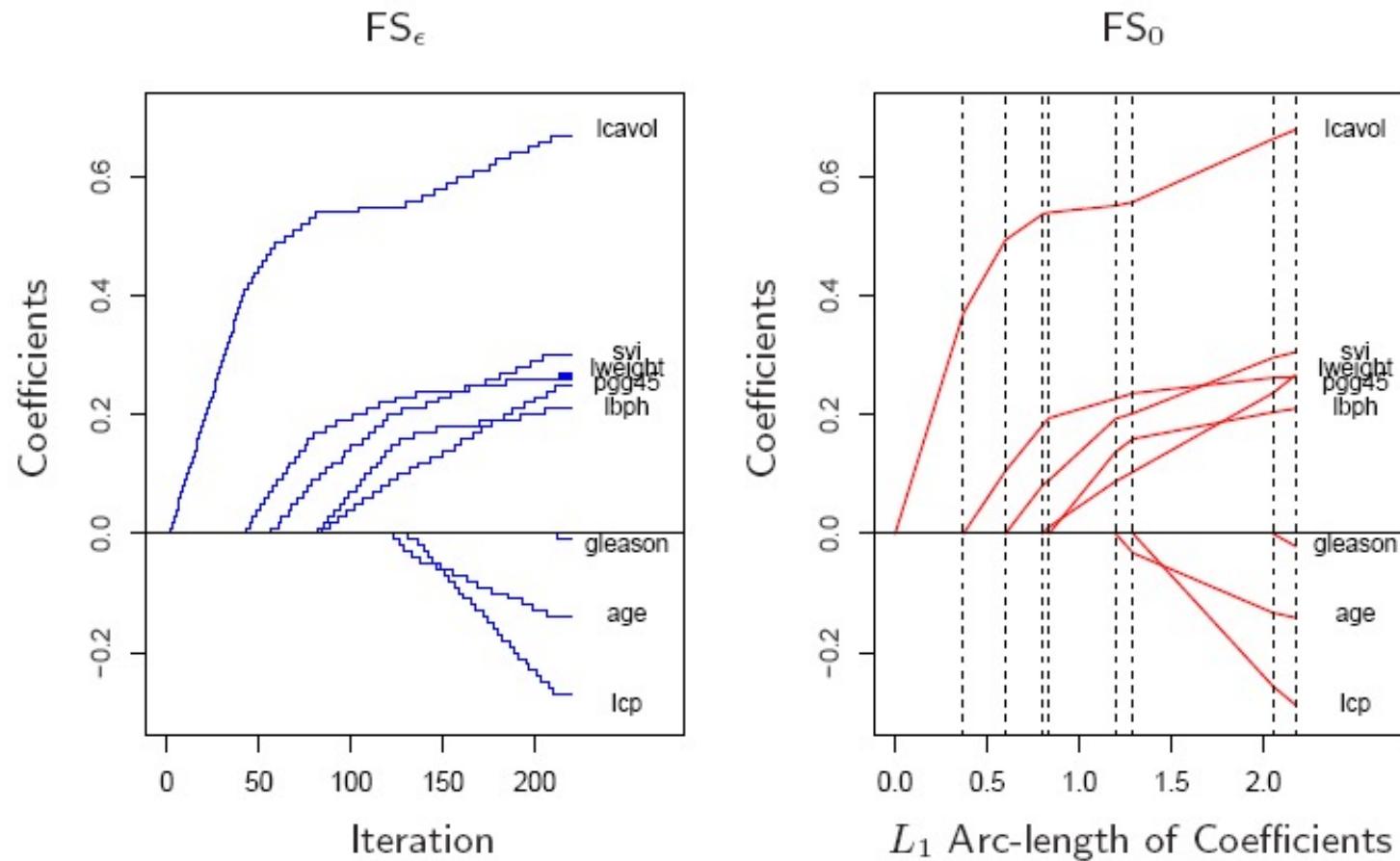


FIGURE 3.19. Coefficient profiles for the prostate data. The left panel shows incremental forward stagewise regression with step size $\epsilon = 0.01$. The right panel shows the infinitesimal version FS_0 obtained letting $\epsilon \rightarrow 0$. This profile was fit by the modification 3.2b to the LAR Algorithm 3.2. In this example the FS_0 profiles are monotone, and hence identical to those of lasso and LAR.

Regression on Derived Input Projections (HTF 3.5)

- PCR: Principal Component Regression
- PLS: Partial Least Squares

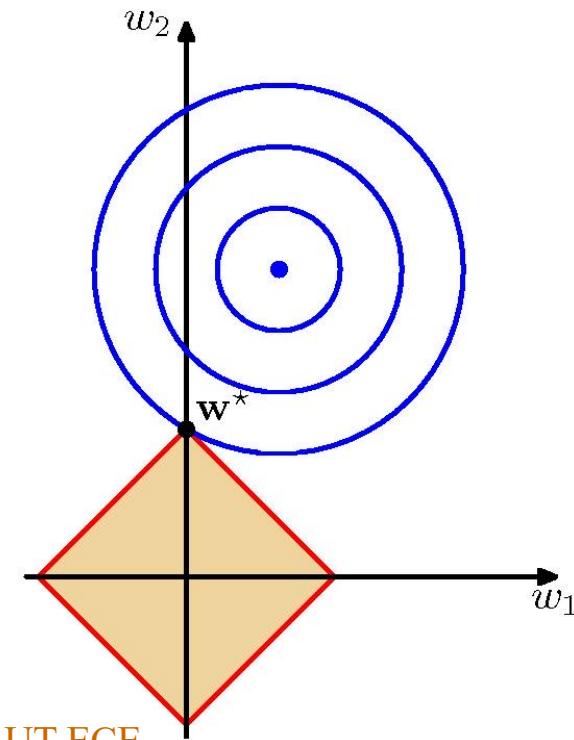
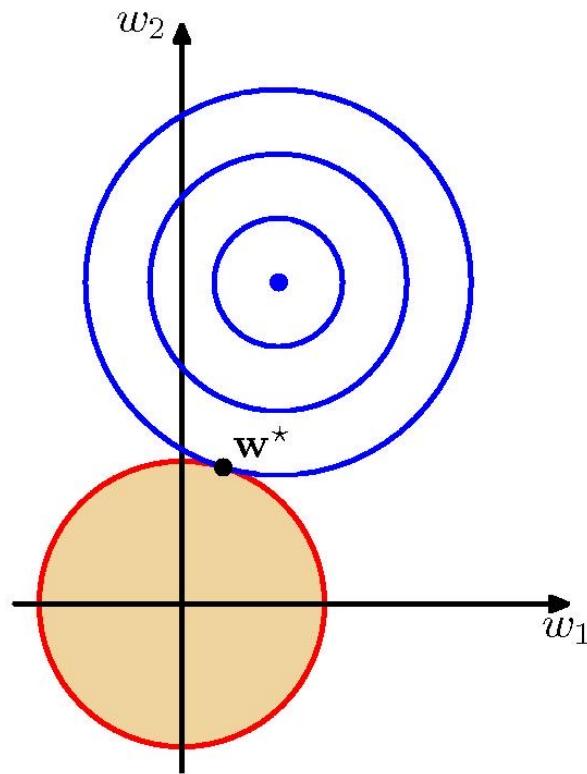
pls package contains both. Typically their effect is similar to ridge regression, but shrinking of coefficients not so smooth

Algorithm 3.3 *Partial Least Squares.*

1. Standardize each \mathbf{x}_j to have mean zero and variance one. Set $\hat{\mathbf{y}}^{(0)} = \bar{y}\mathbf{1}$, and $\mathbf{x}_j^{(0)} = \mathbf{x}_j$, $j = 1, \dots, p$.
2. For $m = 1, 2, \dots, p$
 - (a) $\mathbf{z}_m = \sum_{j=1}^p \hat{\varphi}_{mj} \mathbf{x}_j^{(m-1)}$, where $\hat{\varphi}_{mj} = \langle \mathbf{x}_j^{(m-1)}, \mathbf{y} \rangle$.
 - (b) $\hat{\theta}_m = \langle \mathbf{z}_m, \mathbf{y} \rangle / \langle \mathbf{z}_m, \mathbf{z}_m \rangle$.
 - (c) $\hat{\mathbf{y}}^{(m)} = \hat{\mathbf{y}}^{(m-1)} + \hat{\theta}_m \mathbf{z}_m$.
 - (d) Orthogonalize each $\mathbf{x}_j^{(m-1)}$ with respect to \mathbf{z}_m : $\mathbf{x}_j^{(m)} = \mathbf{x}_j^{(m-1)} - [\langle \mathbf{z}_m, \mathbf{x}_j^{(m-1)} \rangle / \langle \mathbf{z}_m, \mathbf{z}_m \rangle] \mathbf{z}_m$, $j = 1, 2, \dots, p$.
3. Output the sequence of fitted vectors $\{\hat{\mathbf{y}}^{(m)}\}_1^p$. Since the $\{\mathbf{z}_\ell\}_1^m$ are linear in the original \mathbf{x}_j , so is $\hat{\mathbf{y}}^{(m)} = \mathbf{X} \hat{\beta}^{\text{pls}}(m)$. These linear coefficients can be recovered from the sequence of PLS transformations.

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



Estimating True Performance (Formula Driven)*

- true mean squared error ($MSE = SSE/N$) = empirical error + complexity term
 - complexity term = f (model type, # of parameters, # of training points)
 - e.g. linear regression with N samples, P parameters
Akaike's Final Prediction error = $MSE_{\text{empirical}}(N+P) / (N - P)$
 - 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

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 $\mathbf{r} = \mathbf{y} - \bar{\mathbf{y}}$, $\beta_1, \beta_2, \dots, \beta_p = 0$.
 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*