

CS 340 Lec. 15: Linear Regression

AD

February 2011

- Assume you are given some training data $\{\mathbf{x}^i, \mathbf{y}^i\}_{i=1}^N$ where $x^i \in \mathbb{R}^d$ and $\mathbf{y}^i \in \mathbb{R}^c$.
- Given an input test data \mathbf{x} , you want to predict/estimate the output \mathbf{y} associated to \mathbf{x} .
- Applications:
 - \mathbf{x} : location, y sensor reading.
 - \mathbf{x} : stock at time $t - d, t - d + 1, \dots, t - 1$, y : stock at time t .
 - \mathbf{x} : temperature at day $t - d, t - d + 1, \dots, t - 1$, y : temperature at day t .

- We want to learn a mapping/predictor based on $\{\mathbf{x}^i, \mathbf{y}^i\}_{i=1}^N$:

$$f : \mathbb{R}^d \rightarrow \mathbb{R}^c$$

which allows us to predict the response \mathbf{y} given a new input \mathbf{x} ; i.e.

$$\mathbf{y}(\mathbf{x}) = f(\mathbf{x}).$$

- Linear regression is the simplest approach to build such a mapping and is ubiquitous in applied science.

Linear Regression

- For sake of simplicity, consider the simplest case $c = 1$ and $d = 1$.
- Linear regression assumes a model

$$y(x) = w_1x + w_0$$

where $w_1, w_0 \in \mathbb{R}$.

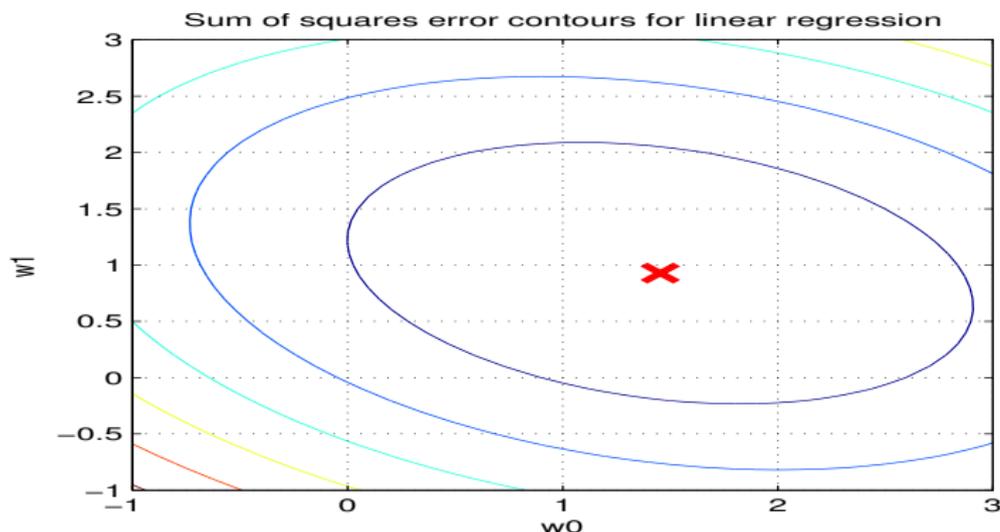
- Given only 2 training data, we can solve for w_1 and w_0 but the result would be dependent of the 2 training data selected and very sensitive to the noise in the training responses y^i .
- A more sensible approach is to minimize the residual errors ε^i over the N training data

$$\begin{aligned}\varepsilon^i &= y^i - y(x^i) \\ &= y^i - w_1x^i - w_0\end{aligned}$$

Least square Regression

- We propose to minimize the sum of the squared residual errors w.r.t (w_0, w_1)

$$E(w_0, w_1) = \sum_{i=1}^N (y^i - w_1 x^i - w_0)^2$$



Least square Regression

- We compute $\frac{\partial E}{\partial w_0}$ and set it equal to zero

$$\begin{aligned}\frac{\partial E}{\partial w_0} &= -2 \sum_{i=1}^N (y^i - w_1 x^i - w_0) = 0 \\ \Leftrightarrow w_0 &= \frac{\sum_{i=1}^N y^i}{N} - w_1 \frac{\sum_{i=1}^N x^i}{N} = \bar{y} - w_1 \bar{x}\end{aligned}$$

- Substituting back $w_0 = \bar{y} - w_1 \bar{x}$ in $E(w_0, w_1)$, we obtain

$$E(\bar{y} - w_1 \bar{x}, w_1) = \sum_{i=1}^N ((y^i - \bar{y}) - w_1 (x^i - \bar{x}))^2$$

Least square Regression

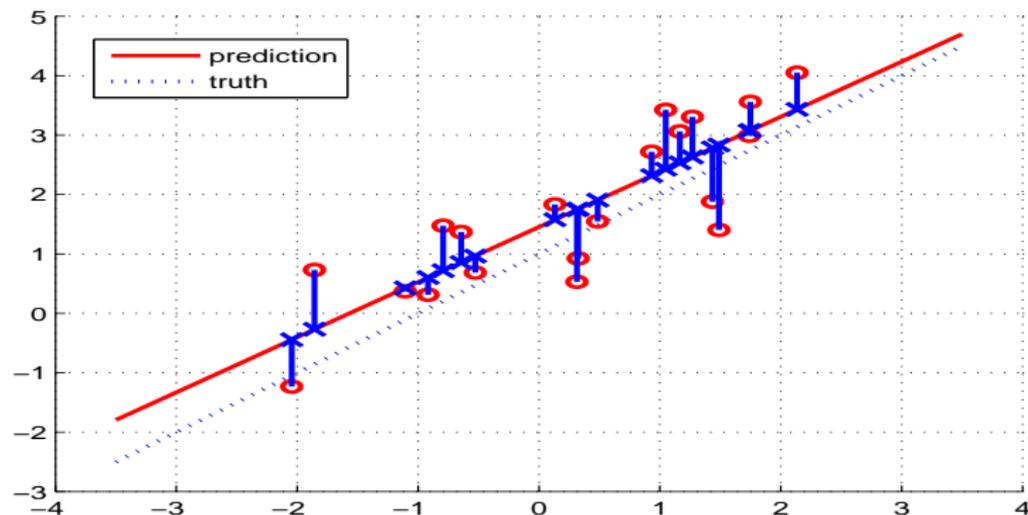
- We compute $\frac{\partial E}{\partial w_1}$ set them equal to zero

$$\begin{aligned}\frac{\partial E}{\partial w_1} &= -2 \sum_{i=1}^N (x^i - \bar{x}) ((y^i - \bar{y}) - w_1 (x^i - \bar{x})) = 0 \\ \Leftrightarrow w_1 &= \frac{\sum_{i=1}^N (x^i - \bar{x}) (y^i - \bar{y})}{\sum_{i=1}^N (x^i - \bar{x})^2}\end{aligned}$$

- Hence

$$w_0 = \bar{y} - w_1 \bar{x} = \bar{y} - \frac{\sum_{i=1}^N (x^i - \bar{x}) (y^i - \bar{y})}{\sum_{i=1}^N (x^i - \bar{x})^2} \bar{x}.$$

Example



In linear least squares, we minimize the sum of squared distances from each training point (denoted by a red circle) to its approximation (denoted by a blue cross). The red diagonal line represents $y(x) = w_1x + w_0$, which is the least squares regression line. Note that these residual lines are not perpendicular to the least squares line, in contrast to PCA.

Least square Regression for Multidimensional Inputs

- Consider now that $c = 1$ but $d > 1$ and we consider

$$\begin{aligned}y(\mathbf{x}) &= w_0 + \sum_{k=1}^d w_k x_k \\ &= \mathbf{w}^T \mathbf{x} + w_0 = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}\end{aligned}$$

where

$$\begin{aligned}\tilde{\mathbf{w}}^T &= (w_0 \ \mathbf{w}^T) = (w_0 \ w_1 \ \cdots \ w_d) \\ \tilde{\mathbf{x}}^T &= (1 \ \mathbf{x}^T) = (1 \ x_1 \ \cdots \ x_d)\end{aligned}$$

- Given N training data, we want to minimize w.r.t $\tilde{\mathbf{w}}$ the sum of the squared residual errors

$$E(\tilde{\mathbf{w}}) = \sum_{i=1}^N (y^i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}^i)^2$$

Least square Regression for Multidimensional Inputs

- We can rewrite

$$E(\tilde{\mathbf{w}}) = (\mathbf{Y} - \tilde{\mathbf{X}}\tilde{\mathbf{w}})^T (\mathbf{Y} - \tilde{\mathbf{X}}\tilde{\mathbf{w}})$$

where

$$\begin{aligned}\mathbf{Y} &= (y^1 \dots y^N)^T && N \times 1 \text{ matrix} \\ \tilde{\mathbf{X}} &= (\tilde{\mathbf{x}}_1 \tilde{\mathbf{x}}_2 \dots \tilde{\mathbf{x}}_N)^T && N \times (d+1) \text{ matrix}\end{aligned}$$

- We can rewrite

$$E(\tilde{\mathbf{w}}) = \mathbf{Y}^T \mathbf{Y} - 2\mathbf{Y}^T \tilde{\mathbf{X}}\tilde{\mathbf{w}} + \tilde{\mathbf{w}}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}}\tilde{\mathbf{w}}$$

- By setting $\frac{\partial E}{\partial \tilde{\mathbf{w}}} = 0$, we obtain if $(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})$ is invertible

$$\hat{\tilde{\mathbf{w}}}_{LS} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{Y}$$

Least square Regression for Multidim Inputs Ouputs

- Consider the case where $c > 1$ and $d > 1$ then

$$\begin{aligned} E(\tilde{\mathbf{w}}) &= \sum_{i=1}^N \text{trace} \left(\mathbf{y}^i - \tilde{\mathbf{W}}^T \tilde{\mathbf{x}}^i \right)^T \left(\mathbf{y}^i - \tilde{\mathbf{W}}^T \tilde{\mathbf{x}}^i \right) \\ &= \left\| \mathbf{Y} - \tilde{\mathbf{X}} \tilde{\mathbf{W}} \right\|_F^2 = \sum_{i=1}^N \sum_{j=1}^c \left(y_j^i - \left(\tilde{\mathbf{W}}^T \tilde{\mathbf{x}}^i \right)_j \right)^2 \end{aligned}$$

where $\tilde{\mathbf{W}} \in \mathbb{R}^{(d+1) \times c}$.

- It can also be shown that

$$\hat{\tilde{\mathbf{w}}}_{LS} = \left(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T \mathbf{Y}$$

Geometric Interpretation

- For sake of simplicity, let's come back to the case where $c = 1$. We have for an new input \mathbf{x}

$$y(\mathbf{x}) = \widehat{\mathbf{w}}_{LS}^T \tilde{\mathbf{x}} = \tilde{\mathbf{x}}^T \widehat{\mathbf{w}}_{LS} = \tilde{\mathbf{x}}^T \left(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T \mathbf{Y}$$

- On the training set $\{\mathbf{x}^i, \mathbf{y}^i\}_{i=1}^N$, we have

$$\mathbf{Y}(\mathbf{X}) = \tilde{\mathbf{X}} \widehat{\mathbf{w}}_{LS} = \tilde{\mathbf{X}} \left(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T \mathbf{Y}$$

so

$$\begin{aligned} \tilde{\mathbf{X}}^T (\mathbf{Y}(\mathbf{X}) - \mathbf{Y}) &= \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \left(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T \mathbf{Y} - \tilde{\mathbf{X}}^T \mathbf{Y} \\ &= 0 \end{aligned}$$

- $\mathbf{Y}(\mathbf{X})$ is simply the orthogonal projection of \mathbf{Y} onto the columns of $\tilde{\mathbf{X}}$.

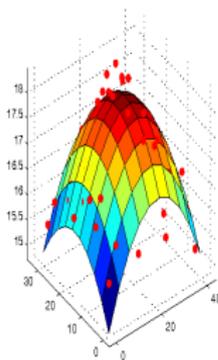
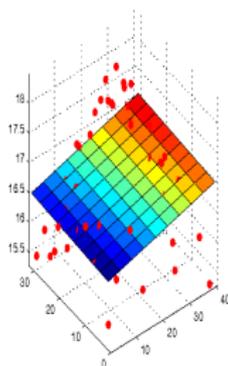
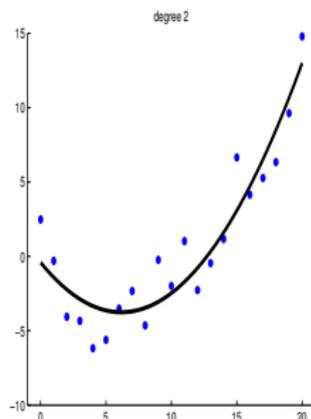
Nonlinear Regression using Basis Functions

- Linear regression can handle nonlinear regression problem; e.g.

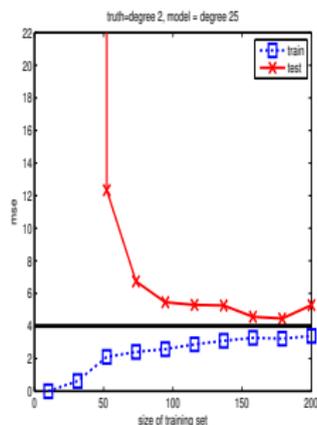
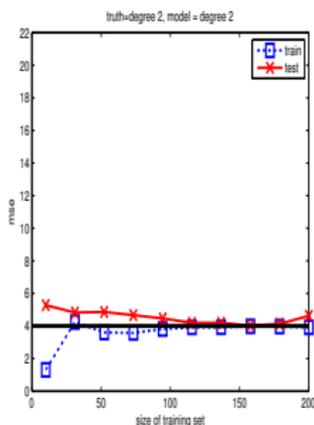
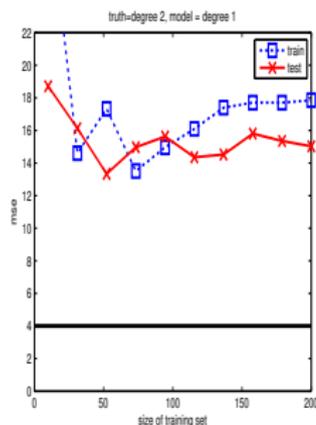
$$y(\mathbf{x}) = \tilde{\mathbf{w}}^T \Phi(\tilde{\mathbf{x}})$$

where $\Phi : \mathbb{R}^{d+1} \rightarrow \mathbb{R}^m$ and $\tilde{\mathbf{w}}$ is a m -dimensional vector.

- Example: $\Phi(\tilde{\mathbf{x}}) = (1, x, x^2)^T$ ($d = 1, m = 3$); $\Phi(\tilde{\mathbf{x}}) = \tilde{\mathbf{x}} = (1, x_1, x_2)$ ($d = 2, m = 3$); $\Phi(\tilde{\mathbf{x}}) = \tilde{\mathbf{x}} = (1, x_1, x_2, x_1^2, x_2^2)$ ($d = 2, m = 5$).



MSE on Training and Test Sets



Data generated from a degree 2 poly. with Gaussian noise of var. 4. We fit polynomial models. For N small, test error of the degree 25 poly. is higher than that of the degree 2 poly., due to overfitting, but this difference vanishes once as N increases. Note also that the degree 1 poly. is too simple and has high test error even given large N .

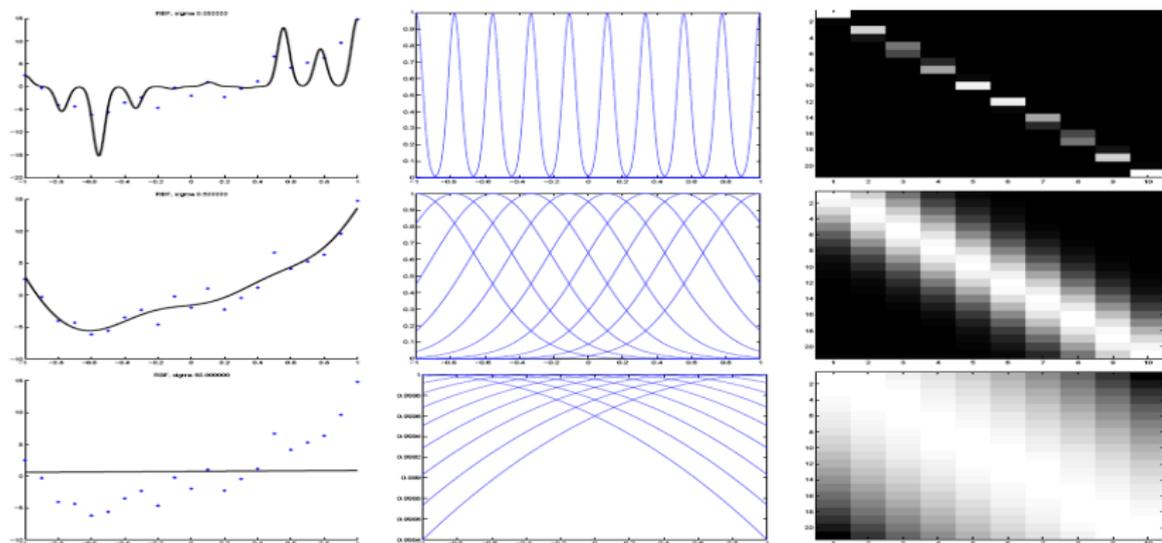
Kernel Regression

- Another way to perform nonlinear regression is to use kernels to define the basis functions $\Phi(\mathbf{x}) = (K(\mathbf{x}, \boldsymbol{\mu}_1), \dots, K(\mathbf{x}, \boldsymbol{\mu}_m))$.
- For example, we could use a Radial Basis Function

$$K(\mathbf{x}, \boldsymbol{\mu}) = \exp\left(-\frac{(\mathbf{x} - \boldsymbol{\mu})^\top (\mathbf{x} - \boldsymbol{\mu})}{2\sigma^2}\right).$$

- Alternatively we can use any function: wavelets, curvelets, splines etc.
- Selecting $(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_m, \sigma^2)$ can be difficult.
- How to select $\boldsymbol{\mu}$: 1) place the centers uniformly spaced in the region containing the data, 2) place one center at each data point, 3) cluster the data and use one center for each cluster, 4) use CV, MLE or Bayesian approach.
- How to select σ^2 : 1) use average squared distances to neighboring centers (scaled by a constant), 2) use CV, MLE or Bayesian approach.

Kernel Regression



(Left) RBS basis in 1d. (Middle) Basis functions evaluated on a grid.
(Right) Design matrix.

A Probabilistic Interpretation of Least-Square Regression

- Consider the following model

$$p(y|\mathbf{x}, \tilde{\mathbf{w}}, \sigma^2) = \mathcal{N}(y; y(\mathbf{x}), \sigma^2)$$

where $y(\mathbf{x}) = \tilde{\mathbf{w}}^T \Phi(\tilde{\mathbf{x}})$.

- Given the training set $\{\mathbf{x}^i, \mathbf{y}^i\}_{i=1}^N$, the MLE estimates of $(\tilde{\mathbf{w}}, \sigma^2)$ is given by

$$\left(\hat{\tilde{\mathbf{w}}}_{MLE}, \hat{\sigma^2}_{MLE} \right) = \arg \max_{(\tilde{\mathbf{w}}, \sigma^2)} \sum_{i=1}^N \log p(y^i | \mathbf{x}^i, \tilde{\mathbf{w}}, \sigma^2)$$

where

$$\sum_{i=1}^N \log p(y^i | \mathbf{x}^i, \tilde{\mathbf{w}}, \sigma^2) = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N \left(y^i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}^i \right)^2$$

- Hence $\hat{\tilde{\mathbf{w}}}_{MLE} = \left(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T \mathbf{Y}$ and $\hat{\sigma^2}_{MLE} = \sum_{i=1}^N \left(y^i - \hat{\tilde{\mathbf{w}}}_{MLE}^T \tilde{\mathbf{x}}^i \right)^2 / N$.

Robust Regression

- The problem with least square regression is that it is very sensitive to outliers as it minimizes

$$E_{L2}(\tilde{\mathbf{w}}) = \sum_{i=1}^N (y^i - y(\mathbf{x}^i))^2 = \sum_{i=1}^N \left(y^i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}^i \right)^2$$

and hence the square of the residual errors.

- To design a procedure less sensitive to these outliers, we could pick

$$E_{L1}(\tilde{\mathbf{w}}) = \sum_{i=1}^N |y^i - y(\mathbf{x}^i)| = \sum_{i=1}^N \left| y^i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}^i \right|.$$

- As $|u|$ goes to infinity slower than u^2 as $|u| \rightarrow \infty$ then this procedure is less sensitive to outliers. We could also use something like

$E_{Huber}(\tilde{\mathbf{w}}) = \sum_{i=1}^N c(y^i - y(\mathbf{x}^i))$ where

$$c(u) = \begin{cases} |u| & \text{if } |u| \leq |u_0| \\ |u_0| & \text{if } |u| \geq |u_0| \end{cases}$$

A Probabilistic Interpretation of Robust Regression

- Consider the following model

$$p(y | \mathbf{x}, \tilde{\mathbf{w}}, b) = \text{Lap}(y; y(\mathbf{x}), b)$$

where $y(\mathbf{x}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$ and

$$\text{Lap}(x; \mu, b) = \frac{1}{2b} \exp\left(-\frac{|x - \mu|}{b}\right)$$

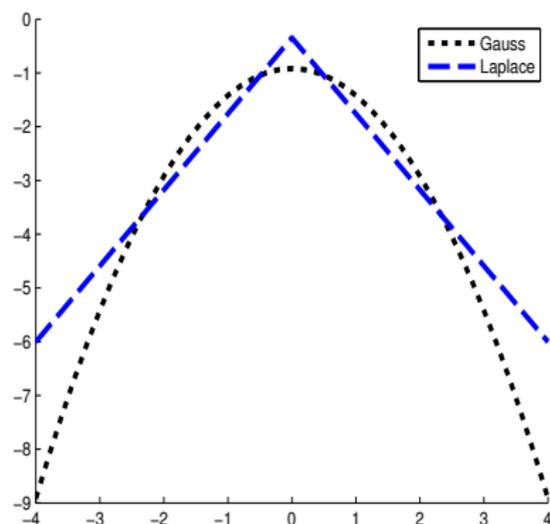
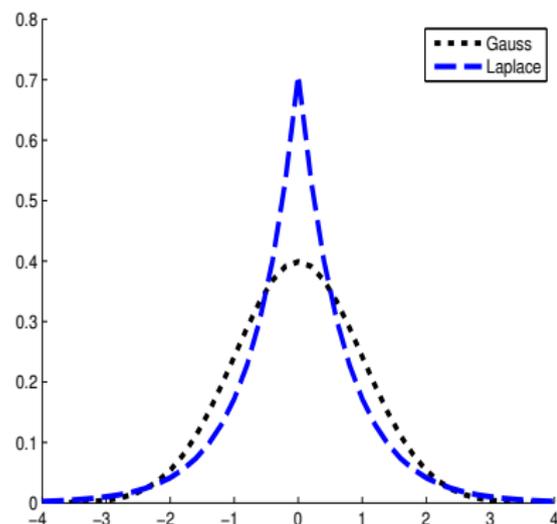
is the Laplace density of location μ and scale b .

- Given the training set $\{\mathbf{x}^i, \mathbf{y}^i\}_{i=1}^N$, the MLE estimates of $(\tilde{\mathbf{w}}, \sigma^2)$ is given by

$$\begin{aligned} \left(\hat{\tilde{\mathbf{w}}}_{MLE, Laplace}, \hat{b}_{MLE} \right) &= \arg \max_{(\tilde{\mathbf{w}}, b)} \sum_{i=1}^N \log p(y^i | \mathbf{x}^i, \tilde{\mathbf{w}}, b) \\ &= \arg \max -N \log(2b) - \frac{1}{b} \sum_{i=1}^N |y^i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}^i|. \end{aligned}$$

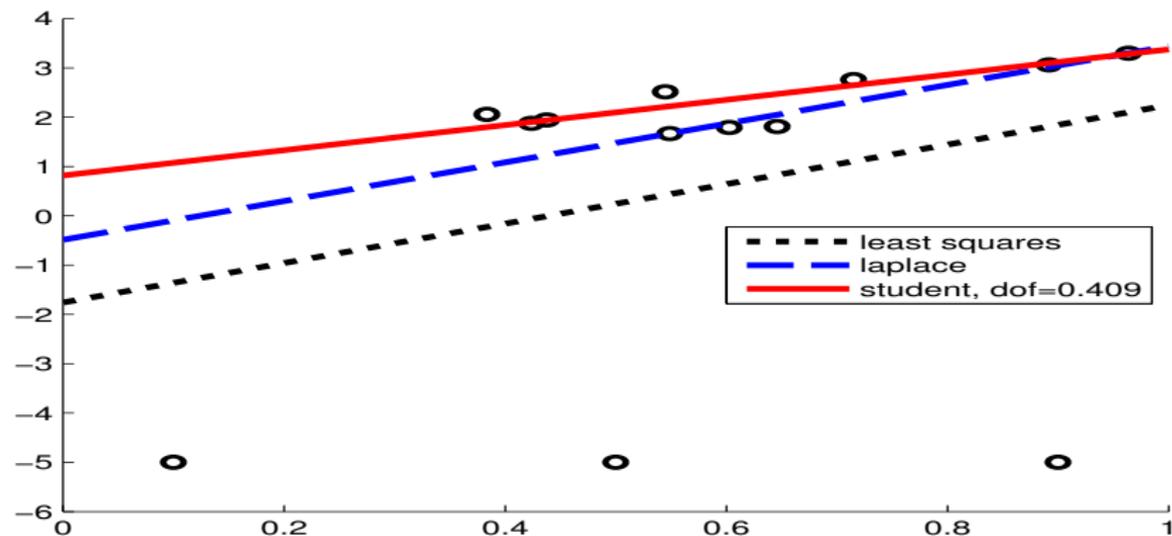
- That is $\hat{\tilde{\mathbf{w}}}_{MLE, Laplace}$ minimizes $E_{L1}(\tilde{\mathbf{w}})$.

Gauss versus Laplace



(left) Plots of $\mathcal{N}(y; 0, 1)$ and $\text{Lap}(y; 0, 1/\sqrt{2})$ densities (both have zero-mean unit variance). (right) Negative logs of these pdfs.

Gauss versus Laplace Regression



Linear regression with outliers. Data points are black circles. Black dotted line: LS, Blue dashed line: Laplace.

Limitations of Least-Square Regression

- Remember that our estimate is given by

$$\hat{\mathbf{w}}_{LS} = \left(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T \mathbf{Y}$$

- This assumes that the $(d + 1) \times (d + 1)$ matrix $\left(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)$ is invertible.
- We have $\text{rank}\left(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}\right) = \text{rank}\left(\tilde{\mathbf{X}}\right) \leq \min(N, d + 1)$. Hence in common scenarios where $N < d + 1$, we can never invert $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$!
- We have also problems when columns/rows of $\tilde{\mathbf{X}}$ are almost linearly dependent (collinearity).

The Need for Regularization

- When we have a small number of data, we want to be able to regularize the solution and limit overfitting.
- When fitting a polynomial regression model, “wiggly” functions will have large weights $\tilde{\mathbf{w}}$.
- For example for the 14 polynomial model fitted previously, we have 11 coefficients w_k such that $|\hat{w}_{k,LS}| > 100!$

Ridge Regression

- To bypass the fact that $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$ might not be invertible, we consider instead

$$\hat{\tilde{\mathbf{w}}}_R = \left(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda I_{d+1} \right)^{-1} \tilde{\mathbf{X}}^T \mathbf{Y}$$

where $\lambda > 0$ and I_{d+1} is the $(d+1)$ identity matrix.

- This estimate minimizes

$$E(\tilde{\mathbf{w}}) = \sum_{i=1}^N \left(y^i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}^i \right)^2 + \lambda \tilde{\mathbf{w}}^T \tilde{\mathbf{w}}.$$

- This is equivalent to minimize

$$\sum_{i=1}^N \left(y^i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}^i \right)^2 \quad \text{s.t.} \quad \tilde{\mathbf{w}}^T \tilde{\mathbf{w}} \leq t(\lambda)$$

- This shrinks the value of $\hat{\tilde{\mathbf{w}}}$ towards zeros.

A Probabilistic Interpretation of Ridge Regression

- Consider the following model

$$p(y | \mathbf{x}, \tilde{\mathbf{w}}, \sigma^2) = \mathcal{N}(y; y(\mathbf{x}), \sigma^2)$$

where $y(\mathbf{x}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$ and

$$p(\tilde{\mathbf{w}} | \sigma^2) = \prod_{k=0}^d \mathcal{N}(w_k; 0, \sigma^2 / \lambda)$$

- In practice, we usually set a flat prior on w_0 ; i.e. $\mathcal{N}(w_0; 0, \sigma^2 / \lambda_0)$ with $\lambda_0 \ll 1$.
- Given $\{\mathbf{x}^i, \mathbf{y}^i\}_{i=1}^N$, the MAP estimate of $\tilde{\mathbf{w}}$ is given by

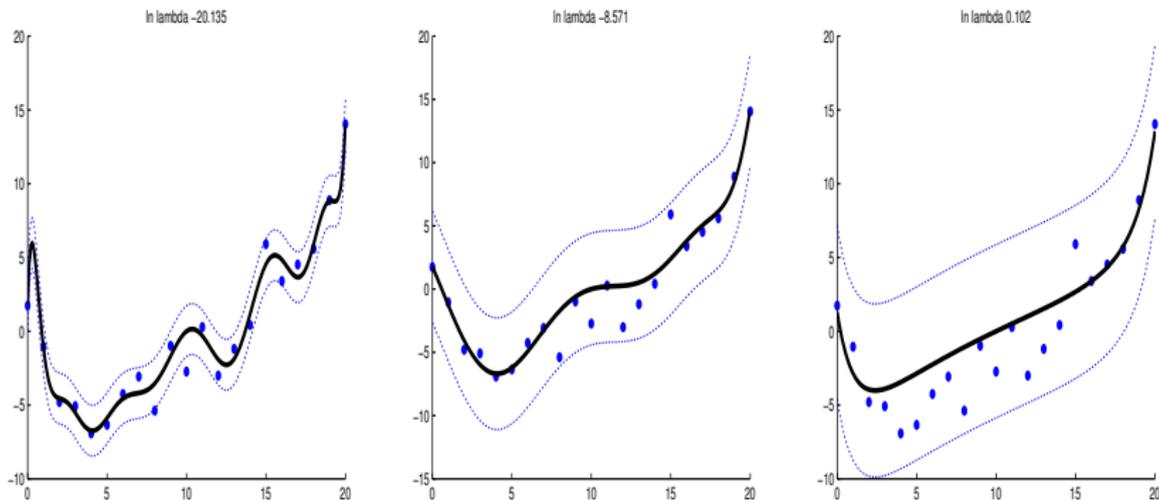
$$\hat{\tilde{\mathbf{w}}}_{MAP} = \arg \max \log p(\tilde{\mathbf{w}} | \mathbf{x}^{1:N}, \mathbf{y}^{1:N}, \sigma^2)$$

where

$$\begin{aligned} \log p(\tilde{\mathbf{w}} | \mathbf{x}^{1:N}, \mathbf{y}^{1:N}, \sigma^2) &= \sum_{i=1}^N \log p(y^i | \mathbf{x}^i, \tilde{\mathbf{w}}, \sigma^2) + \log p(\tilde{\mathbf{w}} | \sigma^2) = \\ &= -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (y^i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}^i)^2 - \frac{\lambda}{2\sigma^2} \tilde{\mathbf{w}}^T \tilde{\mathbf{w}} \end{aligned}$$

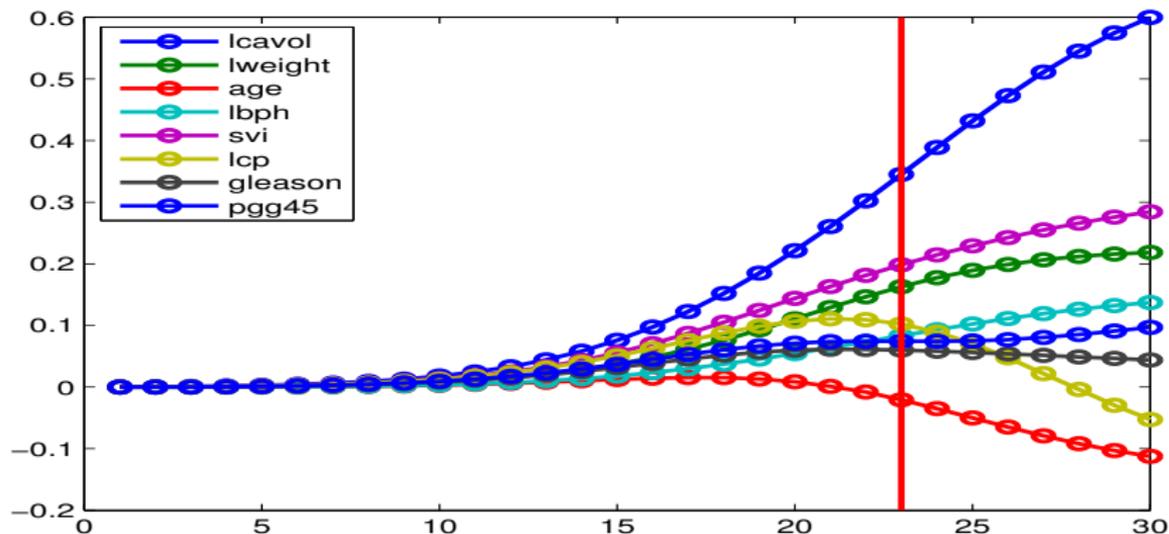
- Hence $\hat{\tilde{\mathbf{w}}}_{MAP} = \hat{\tilde{\mathbf{w}}}_R = \left(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda I_{d+1} \right)^{-1} \tilde{\mathbf{X}}^T \mathbf{Y}$.

Example of Ridge Regression



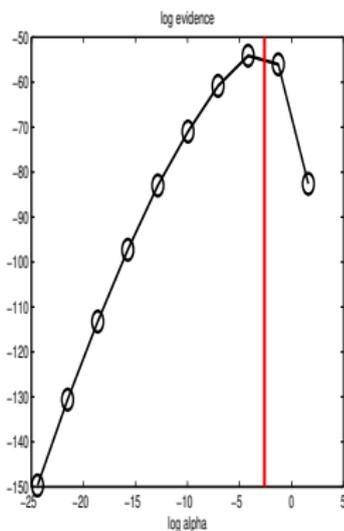
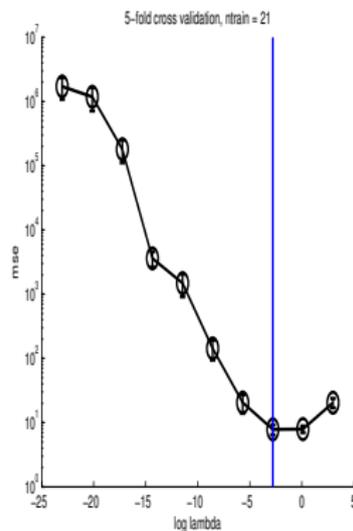
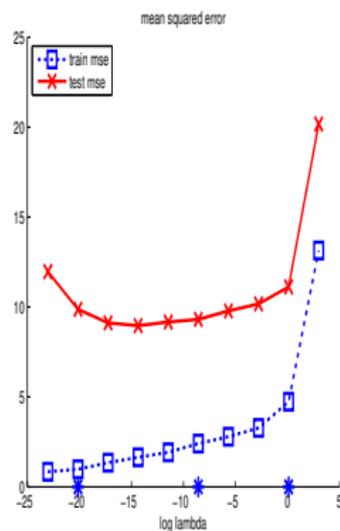
Degree 14 polynomial fit to $N = 21$ data points with increasing λ . The errors bars, represents the standard deviation $\hat{\sigma}$.

Regularization Path for Ridge Regression



Profile of Ridge coefficients for an example on real data where $d = 8$ vs bound on $\mathbf{w}^T \mathbf{w}$, i.e. small $t(\lambda)$ means large λ .

Example of Ridge Regression



(left) Training and test error for a degree 14 poly. with increasing λ .
(center) Estimate of test MSE produced by 5-fold CV. (right)
Log-marginal likelihood vs $\log(\alpha)$ where $\alpha = \lambda/\sigma^2$.

L1 Regression - Lasso

- We minimize in this case

$$E(\tilde{\mathbf{w}}) = \sum_{i=1}^N \left(y^i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}^i \right)^2 + \lambda \|\tilde{\mathbf{w}}\|$$

where $\|\tilde{\mathbf{w}}\| = \sum_{k=0}^d |w_k|$. This is equivalent to minimize

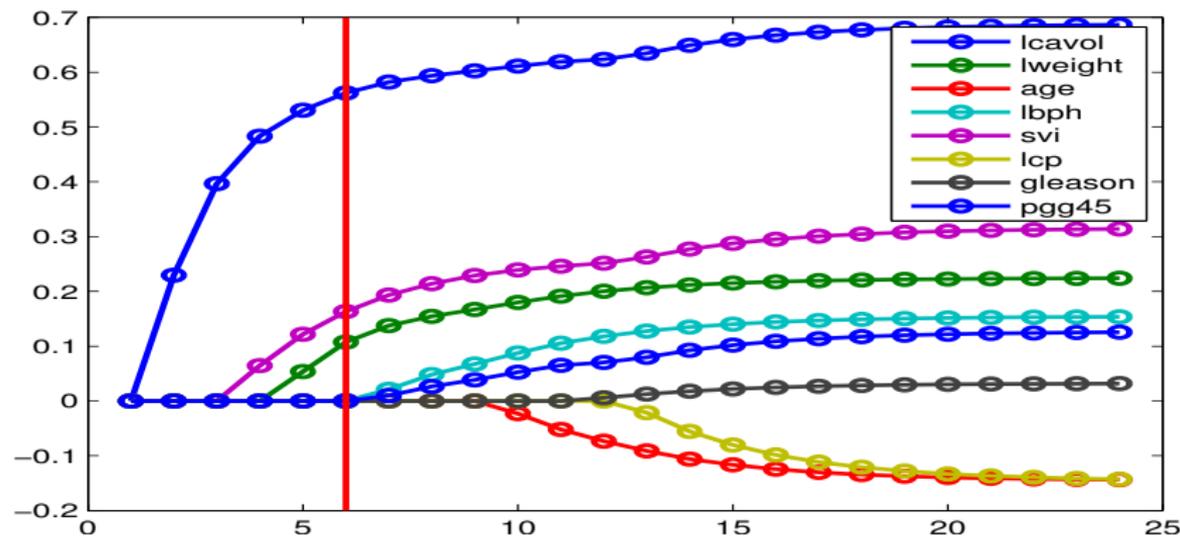
$$\sum_{i=1}^N \left(y^i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}^i \right)^2 \quad \text{s.t.} \quad \|\tilde{\mathbf{w}}\| \leq t(\lambda)$$

- Given $\{\mathbf{x}^i, \mathbf{y}^i\}_{i=1}^N$, this is equivalent to taking the MAP estimate of $\tilde{\mathbf{w}}$ associated to

$$p(y | \mathbf{x}, \tilde{\mathbf{w}}, \sigma^2) = \mathcal{N}(y; y(\mathbf{x}), \sigma^2),$$

$$p(\tilde{\mathbf{w}} | \sigma^2) = \prod_{k=0}^d \text{Lap}(w_k; 0, \sigma^2 / \lambda).$$

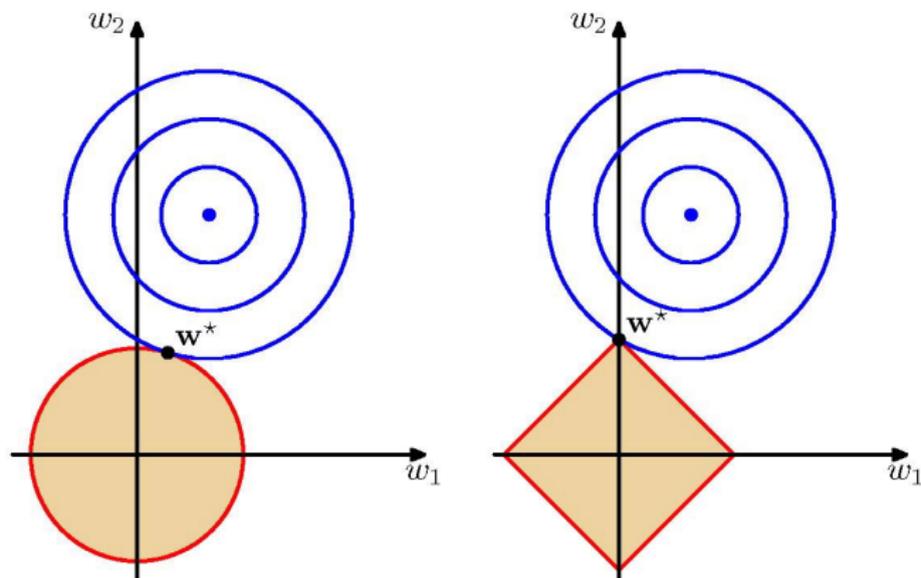
Regularization Path for Lasso



Profile of Lasso coefficients for an example on real data where $d = 8$ vs

bound on $\sum_{k=1}^d |w_k|$, i.e. small $t(\lambda)$ means large λ .

Ridge versus Lasso



Contours of the unregularized function (blue) along with the constraint: ridge (left) and lasso (right). The lasso give a sparse solution as $w_1^* = 0$. The corners of the simplex is more likely to intersect the ellipse than one of the sides as they “stick out” more.