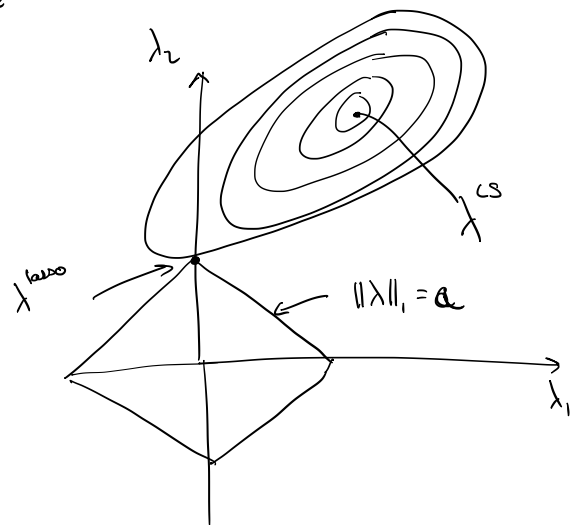
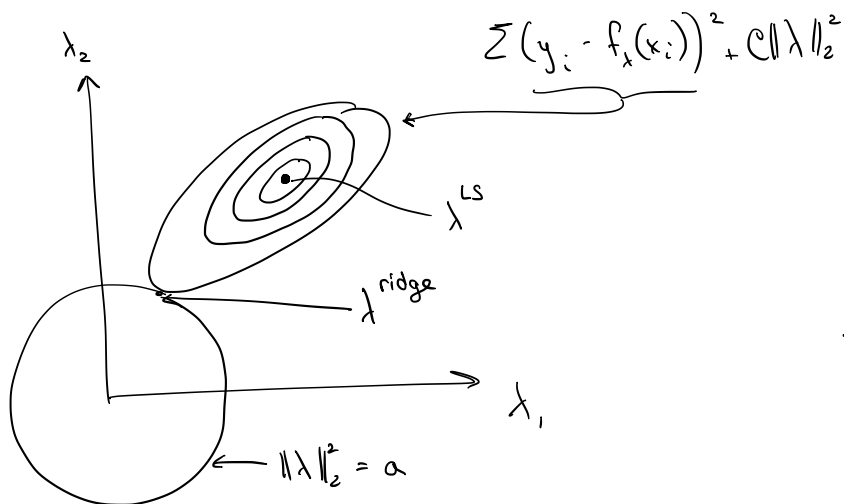


Friends of Ridge Regression

Fact: Lasso does not have a closed form solution,

$$\sum_i (y_i - f_{\lambda}(x_i))^2 + C \|\lambda\|_1$$

Leads to sparser solutions though



## Kernel Least Squares

Regular ridge:

$$F(\bar{\lambda}) = \|\bar{Y} - \bar{X}\bar{\lambda}\|_2^2 + C\|\bar{\lambda}\|_2^2$$

$$\bar{\lambda} \rightsquigarrow \bar{X}^T \bar{d} \quad \leftarrow \text{If I can get } d, \text{ I can get } \lambda$$

$$F(\bar{d}) = \|\bar{Y} - \underbrace{\bar{X}\bar{X}^T}_{n \times n} \bar{d}\|_2^2 + C\|\bar{X}^T \bar{d}\|_2^2$$

$$F(\bar{d}) = \|\bar{Y} - \underbrace{\bar{X}\bar{X}^T}_{K} \bar{d}\|_2^2 + C \underbrace{\langle \bar{X}^T \bar{d}, \bar{X}^T \bar{d} \rangle}_K$$

$$= \|\bar{Y} - K\bar{d}\|_2^2 + C \langle \bar{d}, K\bar{d} \rangle$$

$$F(\bar{d}) = \|\bar{Y} - K\bar{d}\|_2^2 + C \langle \bar{d}, K\bar{d} \rangle$$

$$\nabla F(\bar{d}) = -2K(\bar{Y} - K\bar{d}) + 2C K\bar{d} = 0$$

$$-\bar{Y} + K\bar{d}^* + C\bar{d}^* = 0$$

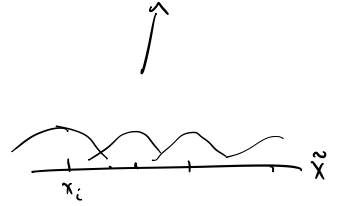
$$\bar{d}^* = (K + C\bar{I})^{-1} \bar{Y}$$

$$X^T = \begin{bmatrix} x_1^T & x_2^T & \dots & x_n^T \end{bmatrix}$$

and remember  $\lambda^* = X^T \bar{d}^* = \sum_i x_i^T \bar{d}_i^*$

Make prediction at new test point  $\tilde{x}$

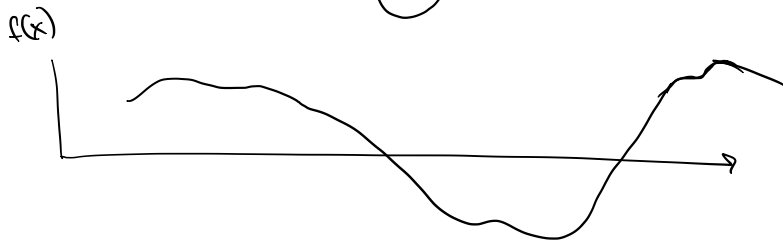
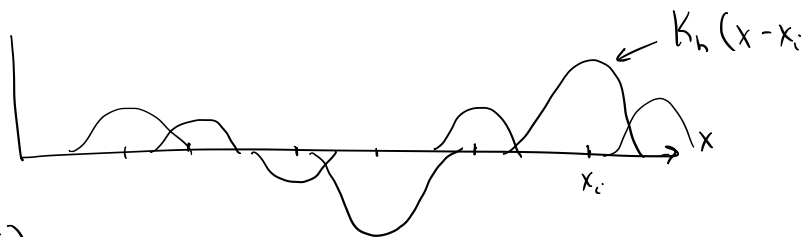
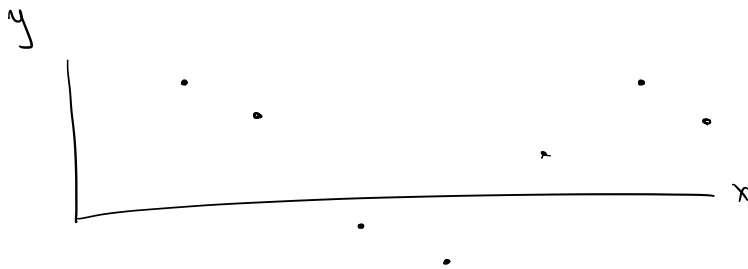
$$f(\tilde{x}) = \tilde{x}^T \lambda^* = \tilde{x}^T \sum_i x_i^T d_i^* = \sum_i \tilde{x}^T x_i^T d_i^* = \sum_i \underbrace{k(\tilde{x}, x_i)}_{\substack{\text{representer theorem} \\ \text{told us this} \\ \text{would happen!!}}} d_i^* = K_{\tilde{x}}^T \cdot d^* = K_{\tilde{x}}^T (K + cI)^{-1} y$$



## Kernel Regression

- There is a separate technique called "kernel regression" that is totally different!  
It is a locally weighted average:

$$f(\tilde{x}) = \frac{\sum_{i=1}^n K_h(\tilde{x} - x_i) y_i}{\sum_{i=1}^n K_h(\tilde{x} - x_i)} \quad \leftarrow \text{Nadaraya-Watson estimator}$$



- Kernel ridge is certainly more sophisticated - actually tries to minimize least square loss

