

Machine Learning I: Supervised Methods

B. Keith Jenkins

Announcements

- Slido event code: 2816950
 - slido.com, Join as participant
- Lecture-time conflicts: sign up
 - See piazza for instructions
- Homework 5 will be posted tomorrow; due on Fri. 3/1.
- Midterm exam in 2 weeks
 - See piazza for materials allowed
- Midterm exam material covered
 - Lectures 1-12, related reading
 - Discussions 1-8, Homeworks 1-5
 - Related piazza posts

Reading

- Bishop 7.1, excluding 7.1.4
 - Support vector machines

Today's lecture

- Regularization
 - Ridge Regression
 - Example
- — — Midterm lecture material ends here — —
- Lagrangian optimization
 - For optimization with constraints
 - Equality constraints
 - Inequality constraint → Lecture 13

Regularization

- Imposes a preference for values of w_j
 - e.g., prefer smaller $|w_j| \forall j$.
 ⇒ modify the criterion function $J(\underline{w})$ to include this preference

Consider l_2 regularization: prefer small $\|\underline{w}\|_2$

Ridge Regression (Regularized Least Squares) [Bishop Sec. 3.1.4]

Add a "regularizer" term $\lambda \|\underline{w}\|_2^2$ to the MSE criterion:

Working in \underline{u} space and \underline{w}' space (dropping primes on \underline{w}):

$$J_{MSE}(\underline{w}') = \frac{1}{N} [\underline{w}^T \underline{u}_i - y_i]^2 = \frac{1}{N} \|\underline{\Phi} \underline{w} - \underline{y}\|_2^2$$

Adding the regularizer term:

$$J_{RR}(\underline{w}') = \frac{1}{N} \|\underline{\Phi} \underline{w} - \underline{y}\|_2^2 + \lambda \|\underline{w}\|_2^2 \quad (\text{augmented notation})$$

This assumes we want to regularize over all components of \underline{w} .

Solve for $\hat{\underline{w}}$:

$$N J_{RR}(\underline{w}) = [\underline{\Phi} \underline{w} - \underline{y}]^T [\underline{\Phi} \underline{w} - \underline{y}] + N \lambda \underline{w}^T \underline{w}$$

$$\begin{aligned} N \nabla_{\underline{w}} J_{RR}(\underline{w}) &= \nabla_{\underline{w}} [\underline{w}^T \underline{\Phi}^T \underline{\Phi} \underline{w} - 2 \underline{w}^T \underline{\Phi}^T \underline{y} + \underline{y}^T \underline{y} + N \lambda \underline{w}^T \underline{w}] \\ &= 2 \underline{\Phi}^T \underline{\Phi} \underline{w} - 2 \underline{\Phi}^T \underline{y} + 2 N \lambda \underline{w} \end{aligned}$$

$$2 \underline{\Phi}^T \underline{\Phi} \hat{\underline{w}} - 2 \underline{\Phi}^T \underline{y} + 2 N \lambda \hat{\underline{w}} = \underline{0}$$

$$(\underline{\Phi}^T \underline{\Phi} + N \lambda \underline{I}) \hat{\underline{w}} = \underline{\Phi}^T \underline{y}$$

$$\begin{aligned} \hat{\underline{w}} &= (\underline{\Phi}^T \underline{\Phi} + N \lambda \underline{I})^{-1} \underline{\Phi}^T \underline{y} \\ \text{Let } \lambda' &= N \lambda \\ \hat{\underline{w}} &= (\underline{\Phi}^T \underline{\Phi} + \lambda' \underline{I})^{-1} \underline{\Phi}^T \underline{y} \end{aligned}$$

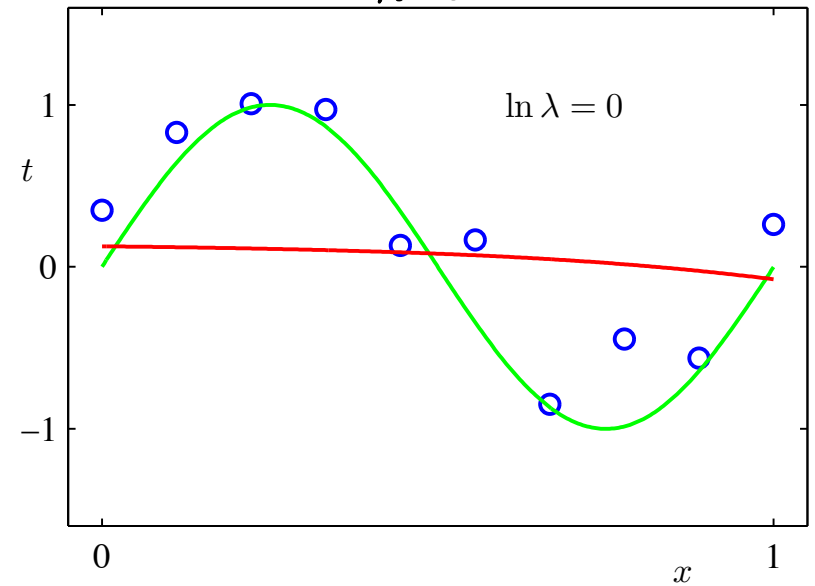
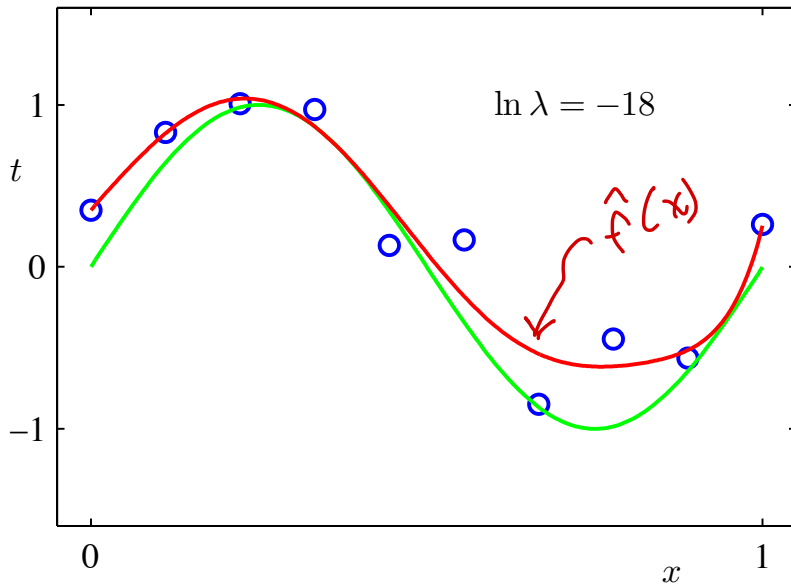
} Ridge
regression
solution

Comment: if choose $\lambda = 0$, then $\hat{\underline{w}}_{RR} = (\underline{\Phi}^T \underline{\Phi})^{-1} \underline{\Phi}^T \underline{y} = \underline{\Phi}^- \underline{y}$

= least-squares regression solution

Regression complexity example: regularization

$N = 10, M = d = 9$ = polynomial order
 $\hat{f}(x) = \sum_{n=0}^9 w'_n x^n$



⑤ If we don't count the regularizer, the above plots have d.o.f. = 10, $N_c = 10$

⑤ Does a regularizer affect d.o.f. or N_c ?

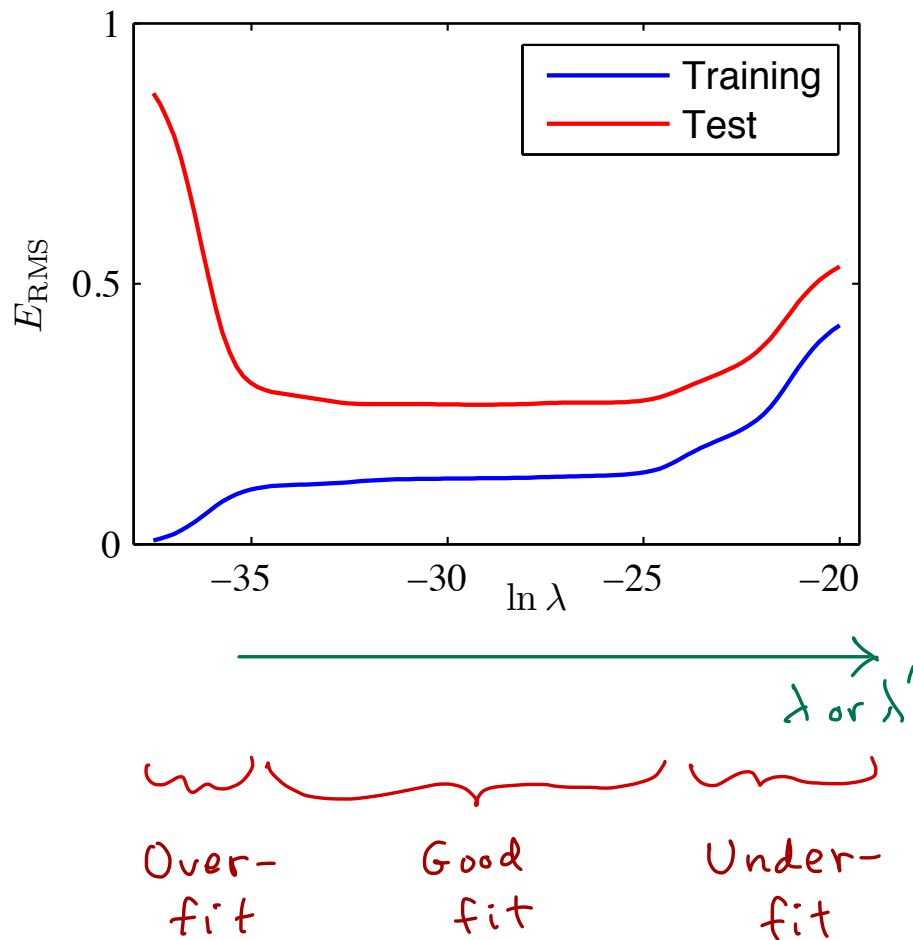
→ Affects N_c because it is adding constraints on the values of w_j .

(similar to data points, which constrain values of w_j by adding eqns $\underline{w}^T \underline{x}_n = b_n$.)

→ Hard to quantify in N_c , so we won't.

Regression complexity example: regularization

$$N = 10, \quad M = d = 9$$



Polynomial Coefficients

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
w_0^*	0.35	0.35	0.13
w_1^*	232.37	4.74	-0.05
w_2^*	-5321.83	-0.77	-0.06
w_3^*	48568.31	-31.97	-0.05
w_4^*	-231639.30	-3.89	-0.03
w_5^*	640042.26	55.28	-0.02
w_6^*	-1061800.52	41.32	-0.01
w_7^*	1042400.18	-45.95	-0.00
w_8^*	-557682.99	-91.53	0.00
w_9^*	125201.43	72.68	0.01

In summary

1. Balancing complexity: d.o.f. vs. N_c
is important for a ML system to perform and generalize well
2. We showed examples of 2 techniques to resolve the case of d.o.f. too large compared with N_c .
 - reduce d.o.f. by reducing dimensionality $D'+1$ of exp. feat. space
 - add a regularizer to prefer some values of w_i (smaller $\| \cdot \|_2^2$) over other values.
 - Called a soft constraint on dimensionality (or on polyn. order).

Next —

Look at another technique to improve on d.o.f. $\leftrightarrow N_c$ balance,
from earlier:

└ Increase N_c by adding other constraints

- Restrict choices of optimal decision boundary [classification]

\Rightarrow Support Vector Machines (SVMs)

Lagrange Optimization with One Equality Constraint [Bishop App.E]

→ Use typical math notation (not ML notation)

Problem: Find an extremum of $f(\underline{x})$ ($=J(\underline{w})$ for us), subject to the constraint:

$$\underbrace{g(\underline{x}) = 0}_{\text{new: constraint}}$$

Solution: 1. Set up a Lagrangian function:

$$L(\underline{x}, \lambda) = f(\underline{x}) + \lambda g(\underline{x})$$

\uparrow original fcn. to minimize (criterion $J(\underline{w})$)
 \uparrow constraint: $g(\underline{x}) = 0$

Note: when constraint is satisfied, $L(\underline{x}, \lambda) = f(\underline{x})$

2. Lagrange method: $\nabla_{\underline{x}, \lambda} L(\underline{x}, \lambda) = \underline{0}$

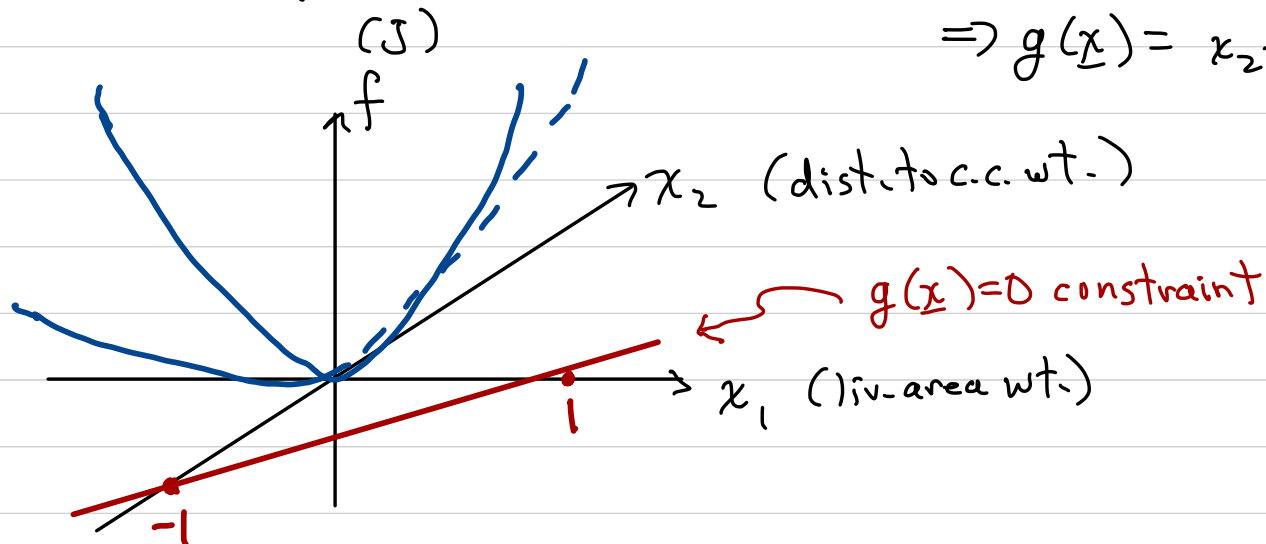
$$\Rightarrow (*) \begin{cases} \nabla_{\underline{x}} L(\underline{x}, \lambda) = \nabla_{\underline{x}} f(\underline{x}) + \nabla_{\underline{x}} \lambda g(\underline{x}) = \underline{0} \\ \nabla_{\lambda} L(\underline{x}, \lambda) = g(\underline{x}) = 0 \end{cases} \leftarrow \text{is our constraint.}$$

[If \underline{x} has dimension d , then $\Rightarrow \begin{cases} d+1 \text{ equations} \\ d+1 \text{ unknowns.} \end{cases}$

Ex. of Lagrangian Optimization with 1 Equality Constraint

Let $f(\underline{x}) = \|\underline{x}\|_2^2 = x_1^2 + x_2^2 \quad (d=2)$

Minimize $f(\underline{x})$ subject to (st.) constraint: $x_2 - x_1 = -1$ ← given constraint
 $\Rightarrow g(\underline{x}) = x_2 - x_1 + 1 = 0$



Lagrange method:

$$\begin{aligned} \text{Let } L(\underline{x}, \lambda) &= f(\underline{x}) + \lambda g(\underline{x}) \\ &= \|\underline{x}\|_2^2 + \lambda (x_2 - x_1 + 1) \end{aligned}$$

$$L(\underline{x}, \lambda) = (x_1^2 + x_2^2) + \lambda (x_2 - x_1 + 1)$$

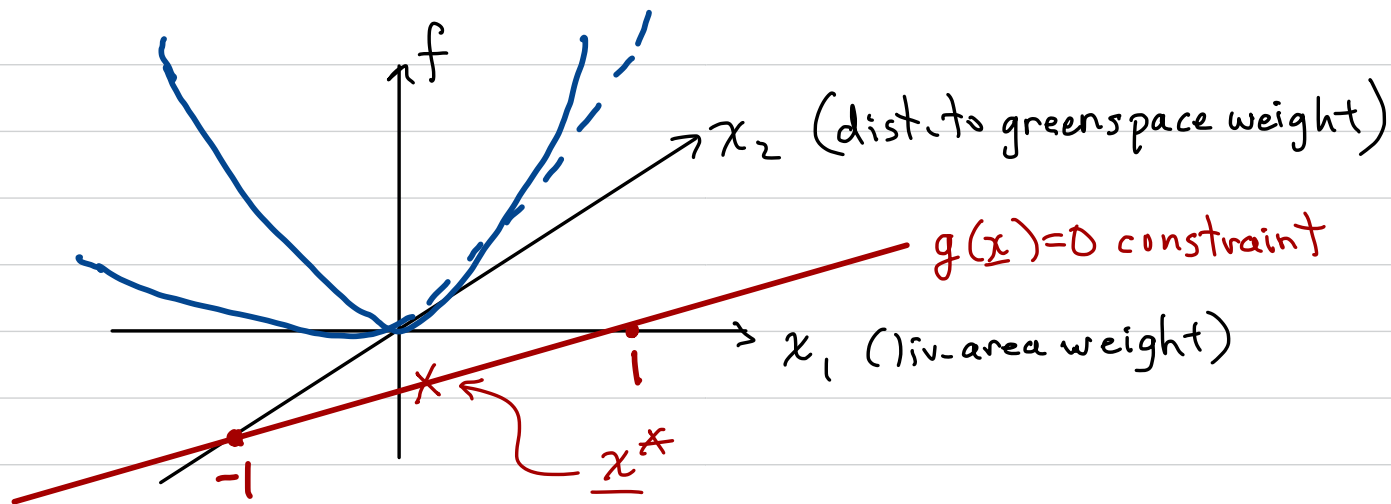
$$\nabla_{\underline{x}, \lambda} L(\underline{x}, \lambda) = \underline{0}$$

$$\begin{cases} \nabla_{\underline{x}} L = \begin{bmatrix} 2x_1 - \lambda \\ 2x_2 + \lambda \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} & \leftarrow (1) \\ \frac{\partial L}{\partial \lambda} = x_2 - x_1 + 1 = 0 & \leftarrow (2) \end{cases}$$

$$\begin{aligned} (1) + (2) &\Rightarrow 2x_1 + 2x_2 = 0 \Rightarrow x_1 = -x_2 \\ \rightarrow (3) &\Rightarrow x_2 - (-x_2) + 1 = 0 \Rightarrow 2x_2 = -1 \Rightarrow \end{aligned}$$

$$\begin{aligned} x_2 &= -\frac{1}{2} \\ x_1 &= +\frac{1}{2} \\ \lambda &= 1 \end{aligned}$$

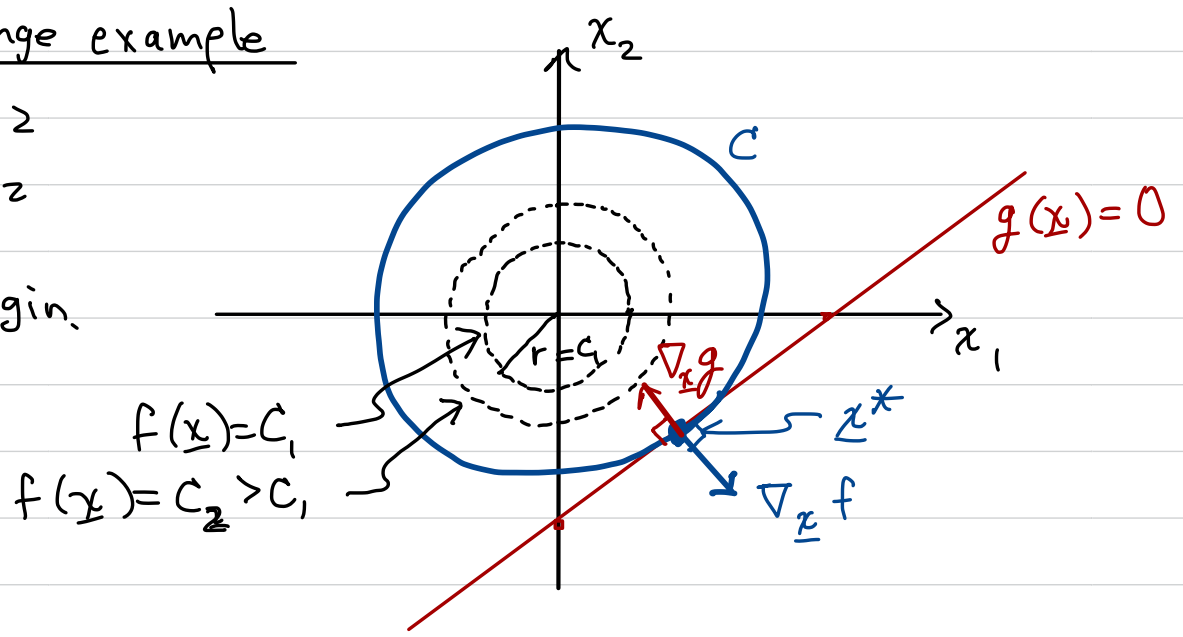
$$\underline{x}^* = \begin{bmatrix} \frac{1}{2} \\ -\frac{1}{2} \end{bmatrix}, \quad \lambda = 1$$



Intuitive explanation of Lagrange example

In x_1, x_2 -plane, $f(\underline{x}) = \|\underline{x}\|_2^2$
 $f(\underline{x}) = d_E^2(\underline{0}, \underline{x}) = r^2$
 $f = \text{Const.}$ is a circle, about origin.

Increase r until circle just touches constraint $g=0$. \rightarrow point \underline{x}^* .
 C is tangent to $g=0$ at this point \underline{x}^* .



$$\left. \begin{aligned} \nabla_{\underline{x}} f(\underline{x}) &\perp f(\underline{x}) = C \\ \nabla_{\underline{x}} g(\underline{x}) &\perp g(\underline{x}) = 0 \\ f(\underline{x}) = C \text{ is tangent to } g(\underline{x}) = 0 \end{aligned} \right\} \text{ at } \underline{x}^*$$

$$\therefore \nabla_{\underline{x}} f(\underline{x}^*) \parallel \nabla_{\underline{x}} g(\underline{x}^*)$$

$$\Rightarrow \boxed{\nabla_{\underline{x}} f(\underline{x}) = \pm \lambda \nabla_{\underline{x}} g(\underline{x}) \text{ at } \underline{x} = \underline{x}^*}$$

$$\nabla_{\underline{x}} f \mp \lambda \nabla_{\underline{x}} g = \underline{0} \quad \text{choose + sign}$$

$$\left\{ \begin{array}{l} \nabla_{\underline{x}} (f + \lambda g) = \underline{0} \\ \nabla_{\underline{x}} (L(\underline{x}, \lambda)) = \underline{0} \\ \text{Coupled with } \nabla_{\lambda} L(\underline{x}, \lambda) = g = 0 \end{array} \right\} \rightarrow \text{gives the opt. sol'n } \underline{x} = \underline{x}^*.$$

Lagrangian Opt'n. with Multiple Equality Constraints

Find min. of $f(\underline{x})$ s.t. $g_i(\underline{x})=0$, $i=1, 2, \dots, R$. $R < d$.

$$L(\underline{x}, \underline{\lambda}) = f(\underline{x}) + \sum_{i=1}^R \lambda_i g_i(\underline{x}) ; \quad \lambda_i = \text{Lagrange multiplier}$$

⑤

$$\nabla_{\underline{x}, \underline{\lambda}} L(\underline{x}, \underline{\lambda}) = \underline{0}$$

$\Rightarrow d+R$ eqns., $d+R$ unknowns