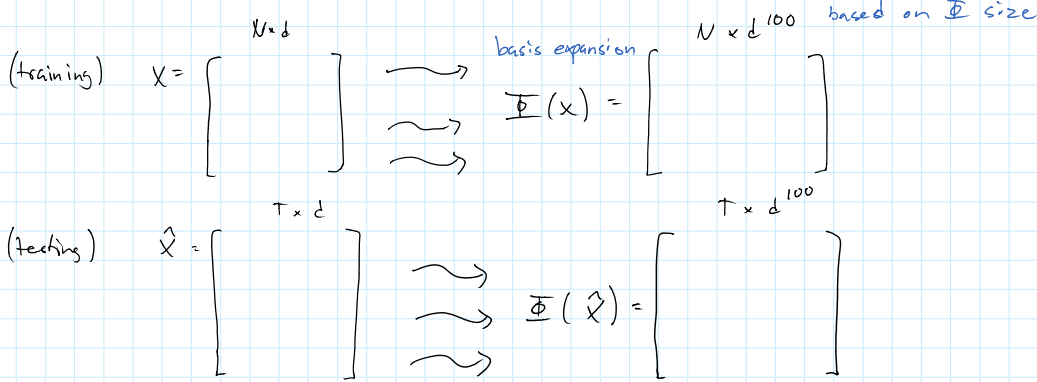


Extra Tutorial:
3:30 → 4:30 Friday (ICCS 193)

Monday, September 22, 2014 3:22 PM

Review Kernel Trick



Ridge regression

$$\begin{aligned} \hat{Y} &= \Phi(\hat{X}) \bar{w}_{MAP} \\ &= \Phi(\hat{X}) (\Phi(X)^T \Phi(X) + \lambda I)^{-1} \Phi(X)^T Y \\ &= \Phi(\hat{X}) \Phi(X)^T (\Phi(X) \Phi(X)^T + \lambda I)^{-1} Y \\ &= k(\hat{X}, X) (K(X, X) + \lambda I)^{-1} Y \end{aligned}$$

$$\bar{w}_{MAP} = (\Phi(X)^T \Phi(X))^{-1} \Phi(X)^T Y$$

matrix inversion lemma (now based on examples, not features)
replace gram with kernel (under conditions)

If you can compute $k(X, X)$ and $k(\hat{X}, X)$,
no need for $\Phi(X)$, $\Phi(\hat{X})$ (don't have to store them)

$$\Phi(X) \Phi(X)^T = \begin{bmatrix} \phi(\bar{x}_1)^T \phi(\bar{x}_1), \phi(\bar{x}_1)^T \phi(\bar{x}_2), \dots \\ \phi(\bar{x}_2)^T \phi(\bar{x}_1), \phi(\bar{x}_2)^T \phi(\bar{x}_2), \dots \\ \vdots \end{bmatrix}$$

$$K(X, X) = \begin{bmatrix} k(\bar{x}_1, \bar{x}_1), k(\bar{x}_1, \bar{x}_2), \dots \\ k(\bar{x}_2, \bar{x}_1), k(\bar{x}_2, \bar{x}_2), \dots \\ \vdots \end{bmatrix}$$

Valid Kernels

mercers theorem

Q: When does there exist feature map ϕ , such that for kernel k we have $\phi(\bar{x}_i)^T \phi(\bar{x}_j) = k(\bar{x}_i, \bar{x}_j)$??

A: If $k(X, X) \succeq 0$ for all X, X_2
(may be hard to show)

positive semi-definite

Checking a valid kernel

Assume k is valid:

- ck $c > 0$
- $k + k'$
- $k(\phi(x), \phi(x'))$
- $\exp(k)$
- $f(k)$ non negative coefficients
- $f(x)k(x, x')f(x')$

polynomial

Course project ideas

- take your application design a kernel
- scale up to large 'N'
- learning a kernel (MKL)

multiple kernel learning

remember for assignment: laplace prior

Search over Feature Combinations

- NP Hard

Greedy Methods:

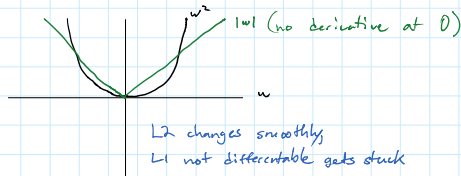
- forward selection: start with empty model, add best next (BIC, cross validation)
(sub-modular: best approximation)
- backwards selection: start with all, remove the worst
- stagewise: combine forward + backwards

L1 - Regularization

$$\underset{w}{\operatorname{argmin}} \frac{1}{2} \|X\bar{w} - Y\|^2 + \lambda \|w\|_1$$

quadratic linear

$$\|w\|_1 = \sum_{i=1}^d |w_i|$$



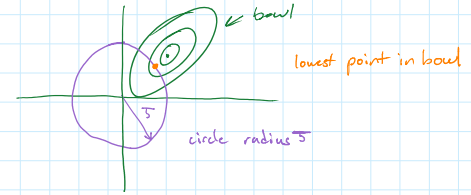
$$\underset{\bar{w}}{\operatorname{argmin}} \frac{1}{2} \|X\bar{w} - Y\|^2$$

s.t. $\|w\|_1 \leq 5$

- regularization: protect against overfitting (but not unique)
- encourages variables to be exactly 0
i.e. not dependant on

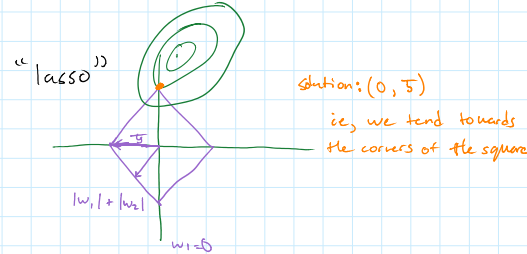
$$\hat{X}w = \begin{bmatrix} \vdots \\ 0 \end{bmatrix}$$

no longer matters



(Smooth) + (non smooth) still easy optimization
but separable

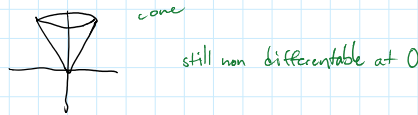
$$f(x) = \sum_{i=1}^d f_i(x_i)$$



Elastic net: $\frac{\lambda_2}{2} \|w\|^2 + \lambda_1 \|w\|_1$

small, for unique unique sparse

Group L1 Regularization: $\lambda \sum_g \|w_g\|$
L1 of groups w_g is a subvector
selects "groups" of variables



$$[\dots \dots \underbrace{0 \ 0 \ 0}_{w_g} \dots \dots]$$

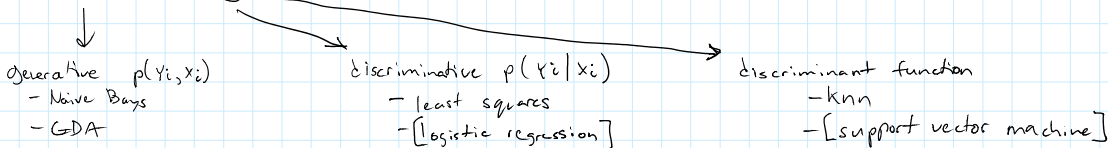
$$X = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \quad Y = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

$N \times d$ $N \times c$ "multiple regression"

$$W = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \quad XW = Y$$

- Nuclear-norm regularization
sum of singular values
 \Rightarrow low rank matrices (many go to zero)
- Structured sparsity ($w_2 \neq 0$ only if $w_1 \neq 0$, etc)
- challenging, but good project (potential publication)
- patterns of variables

Supervised Learning



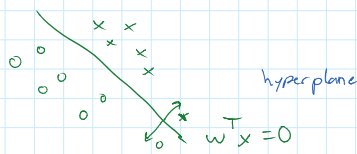
Classification:

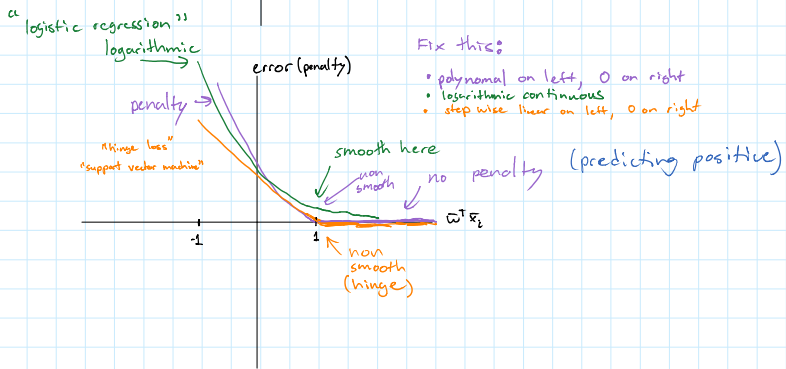
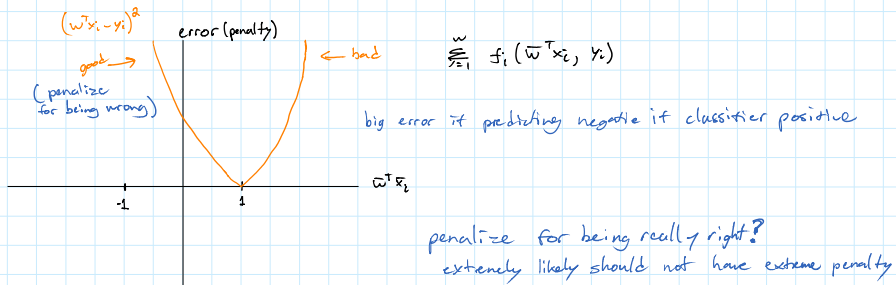
features: $\bar{x}_i \in \mathbb{R}^d$
target: $y_i \in \{-1, 1\}$

Linear classifiers:

$$\hat{y}_i = \operatorname{sign}(w^T \bar{x}_i)$$

sign: $w^T \bar{x}_i > 0 \Rightarrow \hat{y}_i = 1$ (1 if positive)
 $w^T \bar{x}_i < 0 \Rightarrow \hat{y}_i = -1$ (-1 if negative)





Odds ratios

$$\frac{p(\bar{y}_i=1 | \bar{x}_i, \bar{w})}{p(y_i=-1 | \bar{x}_i, \bar{w})}$$

Log odds ratio:

$$\log \left(\frac{p(\bar{y}_i=1 | \bar{x}_i, \bar{w})}{p(y_i=-1 | \bar{x}_i, \bar{w})} \right) = \bar{w}^T \bar{x}_i$$

$$\frac{p(\bar{y}_i=1 | \bar{x}_i, \bar{w})}{\underbrace{p(y_i=-1 | \bar{x}_i, \bar{w})}_{1 - p(y_i=1 | \bar{x}_i, \bar{w})}} = \exp(\bar{w}^T \bar{x}_i)$$

$$p_1 = (1 - p_1) \exp(\bar{w}^T \bar{x}_i)$$

$$p_1 (1 + \exp(\bar{w}^T \bar{x}_i)) = \exp(\bar{w}^T \bar{x}_i)$$

$$p_1 = \frac{\exp(\bar{w}^T \bar{x}_i)}{1 + \exp(\bar{w}^T \bar{x}_i)} = \frac{1}{1 + \exp(-\bar{w}^T \bar{x}_i)}$$

"sigmoid function"

$$p_2 = \frac{1}{1 + \exp(\bar{w}^T \bar{x}_i)}$$