

Converting to constrained problem

- non-smooth optimization is hard
- smooth constrained optimization often simpler

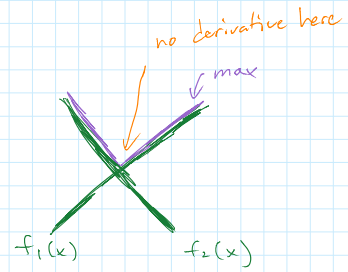
Trick to convert (1 or many):

If we have $\min_x \max_x \{f_i(x)\}$, $\circ \circ \circ$ $\max \{x, -x\} = |x|$

↓ equivalent

$\min_{x,v} v$, subject to $v \geq \max_i \{f_i(x)\}$ (one value)

$v \geq f_i(x), \forall i$ (bigger than functions individually)



Change of Basis

Original data

$$x = \begin{bmatrix} 1 \\ 0.5 \\ 2 \\ \vdots \\ 10 \end{bmatrix}$$

make new X
"Φ(X)"
ϕ: x̄_i → φ(x̄_i)

template
[1 x x²]

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 0.5 & 0.25 \\ 1 & 2 & 4 \\ \vdots & \vdots & \vdots \\ 1 & 10 & 100 \end{bmatrix}$$

new design matrix (then business as usual)

Cross Validation

- large 'k': less chance of overfitting test error, but higher variance + running time
- idea of repeating CV with different random partitioning - lower variance

Ridge Regression

$$f(w) = \frac{1}{2} \|X\bar{w} - Y\|^2 + \frac{\lambda}{2} \|\bar{w}\|^2$$

log likelihood prior

- what? shrink w_i towards 0
- why? - solution is unique
- avoid overfitting due to huge values
- improving conditioning
- "Magic": do a huge basis expansion (lets you use big design matrix without overfitting)

$$\nabla f(w) = X^T X \bar{w} - X^T Y + \lambda \bar{w}$$

(set) = 0 (at the stationary point)

$$(X^T X + \lambda I) \bar{w} = X^T Y$$

identity matrix to align dimensions

$$\bar{w} = (X^T X + \lambda I)^{-1} X^T Y \iff \bar{w} = X^T (X X^T + \lambda I)^{-1} Y$$

matrix inversion lemma

$$\nabla^2 f(w) = X^T X + \lambda I \succ 0$$

positive definite solution is unique

- X^TX: "scatter matrix"
- 1/N X^TX: "covariance MLE"
- X X^T: "Gram matrix"

MAP estimation

MLE: $\operatorname{argmax}_{h \in H} p(D|h)$

MAP: $\operatorname{argmax}_{h \in H} p(h|D)$

hypothesis "is model h reasonable?"

posterior \propto likelihood prior $p(D|h)p(h)$

↑ proportional to. $f(x) \propto x = cx$

$\sum p(x) = 1$ $g(x) \propto \frac{y^2 + \alpha^2}{\sigma^2} x$ still true

"exists constants not depend h," subject to $\sum_h p(h|D) = 1$

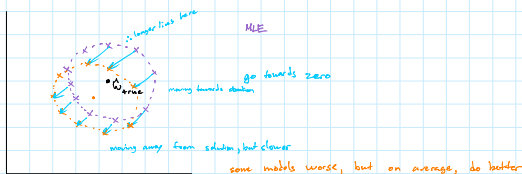
$$\log p(h|D) = \log p(D|h) + \log p(h) + \text{const}$$

$$-\log p(h|D) = -\log p(D|h) - \log p(h) + \text{const}$$

$$= \sum_{i=1}^N (y_i - \bar{x}_i^T \bar{w})^2 + \frac{1}{2\sigma^2} \bar{w}^T \bar{w}$$

"precision"
 $\lambda = \frac{1}{\sigma^2}$ $\frac{\lambda}{2} \|\bar{w}\|^2$

Seems
para-
dox



assumptions: gaussian prior

$$y_i | \bar{x}_i, \bar{w} \sim \mathcal{N}(\bar{w}^T \bar{x}_i, \sigma^2)$$

$$w_i | \sigma^2 \sim \mathcal{N}(\emptyset, \sigma^2)$$

$$p(w_i) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{1}{2} \frac{(w_i - \emptyset)^2}{\sigma^2}\right)$$

$$\sum_{i=1}^n w_i^2 = \|\bar{w}\|^2$$

given σ^2

$$\log p(w_i) = -\frac{1}{2} \frac{(w_i)^2}{\sigma^2} - \underbrace{\log(\sigma) - \log(\sqrt{2\pi})}_{\text{constant}}$$

Kernels

We assume $\bar{x}_i \in \mathbb{R}^d$
 what if we don't know representation?

Ex: "The cat is back"
 "The black cat is back"

Longest common subsequence ("cat is back")
 Edit distance ("cat is back")

"kernel function"

$$k(\bar{x}_i, \bar{x}_j) \mapsto \mathbb{R} \quad \text{often a similarity}$$

$x \in \mathcal{X}$
 some set of all abstract space \mathcal{X}
 "the cat is back" $\in \mathcal{X}$

Typically, $k(\bar{x}_i, \bar{x}_j) \geq 0$
 $k(\bar{x}_i, \bar{x}_j) = k(\bar{x}_j, \bar{x}_i)$

Examples

Linear

"Linear" $\bar{x}_i^T \bar{x}_j$

"Poly" $(\bar{x}_i^T \bar{x}_j + a)^m$

"RBF" $\exp\left(-\frac{\|\bar{x}_i - \bar{x}_j\|^2}{2\sigma^2}\right)$

section 14.2

- string kernels

- pyramid match kernel

How to use?

$$\phi(\bar{x}_i) = [k(\bar{x}_i, \bar{x}_1) \quad k(\bar{x}_i, \bar{x}_2) \quad \dots \quad k(\bar{x}_i, \bar{x}_n)]$$

$$[k(\bar{x}_1, z_1) \quad \dots \quad k(\bar{x}_1, z_m)] \quad \text{man? possible source project}$$

$$\Phi(X) = \begin{bmatrix} k(\bar{x}_1, \bar{x}_1) & \dots \\ \vdots & \\ k(\bar{x}_n, \bar{x}_n) \end{bmatrix} \Rightarrow \text{Gram Matrix}$$

linear: $\Phi(x) = X X^T$

Kernel Trick

Some kernels have an explicit feature map (inner product)

$$k(x, z) = (X^T z)^2$$

$$= (x_1 z_1 + x_2 z_2)^2$$

$$= x_1^2 z_1^2 + 2x_1 x_2 z_1 z_2 + x_2^2 z_2^2$$

$$= (x_1^2, \sqrt{2} x_1 x_2, x_2^2)^T (z_1^2, \sqrt{2} z_1 z_2, z_2^2)$$

$$= \phi(x)^T \phi(z) \quad \text{Kernels between data points}$$

$$\bar{w}_{\text{MAP}} = X^T (X X^T + \lambda I)^{-1} Y$$

$$f(X \bar{w}) + \frac{\lambda \|\bar{w}\|^2}{2}$$

anything written like this,
 can be kernelized

Test time

$$\begin{aligned}\hat{y} &= \hat{X} \bar{w}_{MAP} \\ &= \hat{X} X^T (X X^T + \lambda I)^{-1} Y \\ &\quad k(\hat{X}, X) \quad k(X, X) \quad \text{kernel is } n \text{ by } n \text{ matrix} \\ &\quad k(\hat{X}, X) (k(X, X) + \lambda I)^{-1} Y\end{aligned}$$

"Gaussian Process"

Feature Selection

What if we don't know which features x are relevant

$$x = \begin{bmatrix} \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{bmatrix}$$

Search over features

- NP-hard
- 16, ok. 1million? 😞