

Sparse regression

Tuesday, February 27, 2024 9:05 AM

Today sparse regression and dictionary learning

Problem 1: Given time series data, can we fit dynamics to it in a Bayesian context?

Given $\tilde{y}_d = (y_{t_1}^d, \dots, y_{t_N}^d)$

Hypothesize $\begin{cases} \dot{y} = F(t, y | \theta) \\ y(t=0) = y_0 \end{cases}$

Minimize a norm $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \mathcal{L}(y - (y_0 + \int_0^t F dt))$

Problem 2 Assume direct measurements of relevant DOFs z are unavailable

$$z = G(y | \theta_1)$$

$$\begin{cases} \dot{z} = F(z, t | \theta_2) \\ z(t=0) = G(y(t=0)) \end{cases}$$

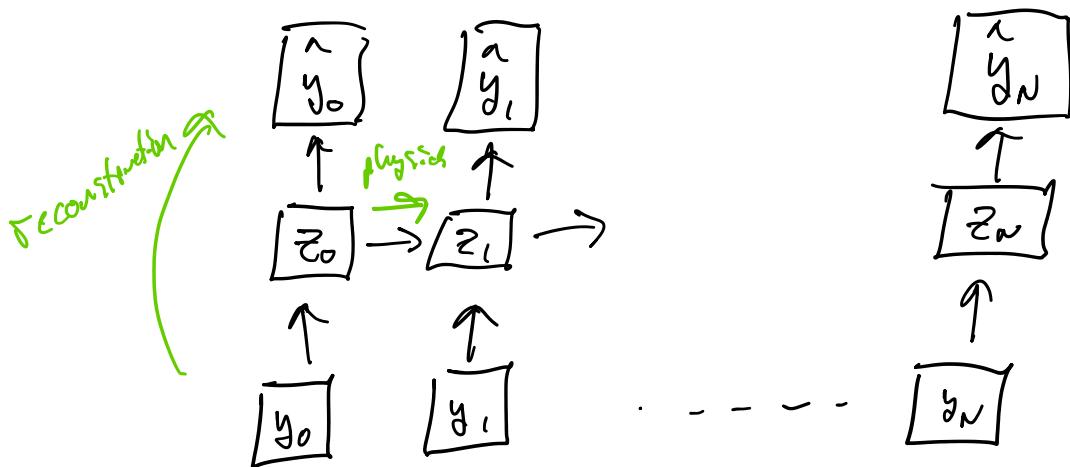
example Image-based simulation

Candidate architectures

Coordinate

$G \rightarrow$ auto encoders

$F \rightarrow$ neural ODEs, universal differential
eqns, dictionary learning



We'll start w/ one simple strategy
called SINDy

"Sparse Identification of nonlinear
dynamics" 2016 Bruckner
Proctor
Kutz

IDEA $\dot{x} = \Theta(x) \xi$

dictionary $\Theta(x)$ ξ unknown coefficients

SPARSITY

Want to only choose a few terms, so data reveals governing ODE
 "Occam's Razor", Parsimony

Sparse Reg vs Polynomial Fitting

Not combining basis functions to approximate more complicated ones

Ex images of pendulum swinging
 Know $\ddot{z} = -kz$ is model

Can we infer θ coordinate from images?

Necessary Background

Revisiting overfitting (Murphy § 4.5, § 11.3)

$$L(\theta; \lambda) = \left[\frac{1}{N} \sum_{i=1}^N l(y_i; \theta; x_i) \right] + \lambda C(\theta)$$

↑ ↑
 Reg. term complex

Typical Bayesian setting

$$l = \log p(y_i | x_i; \theta)$$

$$C(\theta) = \log p(\theta) \leftarrow \text{prior}$$

def Maximum a posteriori (MAP) estimation

$$\hat{\theta} = \arg \max_{\theta} \log p(\theta | \theta)$$

$$\hat{\theta} = \arg \max_{\theta} \log \left(\frac{p(\theta | \theta) p(\theta)}{p(\theta)} \right)$$

$$\hat{\theta} = \arg \max_{\theta} \left[\log p(\theta | \theta) + \log p(\theta) \right]$$

↑ Likelihood ↑ Prior

note the unknown marginal likelihood $p(\theta)$
drops out

Bayesian regression revisited

$$\text{Approximate } g(x) = \sum_i \Phi_i(x) c_i \\ = P(x)^T C$$

$$\text{Model } p(y|x, c) = N(y; \mu = P(x)^T C, \sigma^2)$$

$$\mu_{MLE} = P(x) (P(x)P(x)^T)^{-1} P(x) y$$

Ridge Regression

Adopt the prior

$$p(c) = N(c; 0, \lambda^2 I)$$

$$\begin{aligned} \mathcal{L} &= \log p(y|x, c) p(c) = \\ &\log N(y; P(x)^T C, \sigma^2) N(C; 0, \lambda^2 I) \end{aligned}$$

$$= -\frac{1}{2} \ln \pi \sigma^2 - \frac{1}{2} \underbrace{(P(x)^T C - y)^2}_{+ \lambda^2 C^2}$$

$$= \sum_d -\frac{1}{2} \log \sigma^2 - \frac{1}{2} \left(\frac{\mathbf{P}(x_d)^T \mathbf{C} - y_d}{\sigma^2} \right)^2 - \frac{1}{2} \frac{\mathbf{C}^2}{\lambda^2}$$

$$0 = \nabla_{\mathbf{C}} L = \nabla_{\mathbf{C}} \sum_d \left(\frac{\|\mathbf{P}(x_d)^T \mathbf{C} - y_d\|^2}{\sigma^2} + \frac{1}{\lambda^2} \|\mathbf{C}\|^2 \right)$$

$$0 = \sum_d 2 \left(\mathbf{P}(x_d)^T \mathbf{C} - y_d \right) \mathbf{P}(x_d) + 2 \left(\frac{\sigma^2}{\lambda^2} \mathbf{C} \right)$$

$$\sum_d \left(\mathbf{P}(x_d) \mathbf{P}(x_d)^T + \frac{\sigma^2}{\lambda^2} \mathbf{I} \right) \mathbf{C}_{\text{MAP}} = \sum_d \mathbf{P}(x_d) y_d$$

↑ ridge stabilization

- penalizing magnitude of coefficients

- aka - ridge regression

- L_2 regularization

- weight decay

- Gives invertible matrix in small data limits

For further reading see Murphy §11.7

- Conjugate priors

- Bayesian Linear Regression

Lasso / L' regression for sparsity

- Many applications want coefficients to not just be small but as few nonzero as possible

not just be small but as few nonzero entries as possible

Goal

$$\ell_0 \text{ norm} \quad \|w\|_0 = \sum_{d=1}^D \prod_{|w_d| > 0} \quad \text{minimization}$$

- Important for feature selection & dictionary learning

Lasso Absolute shrinkage + selection operator
(Tibshirani: 96)

$$L(w) = -\log p(y|w) - \log(w/\lambda)$$

if $w|\lambda \sim \prod_d \text{lap}(w_d | \mu=0, b=\lambda^{-1})$

$$\text{lap}(w|\mu, b) = \frac{1}{2b} \exp\left(-\frac{|w-\mu|}{b}\right)$$

$$\text{Then } L(w) = \|Xw - y\|^2 + \lambda \|w\|_1$$

Note that for other norms $\|w\|_p = \left(\sum |w_i|^p\right)^{1/p}$
we get sparser solutions, but $p < 1$ is nonconvex

- ℓ_1 -norm is tightest convex relaxation of ℓ_0 -norm

Why sparse?

$$\min_w -\log p(y|w) + \lambda \|w\|_1$$

is Lagrangian associated w/

$$\min_w -\log p(y|w)$$

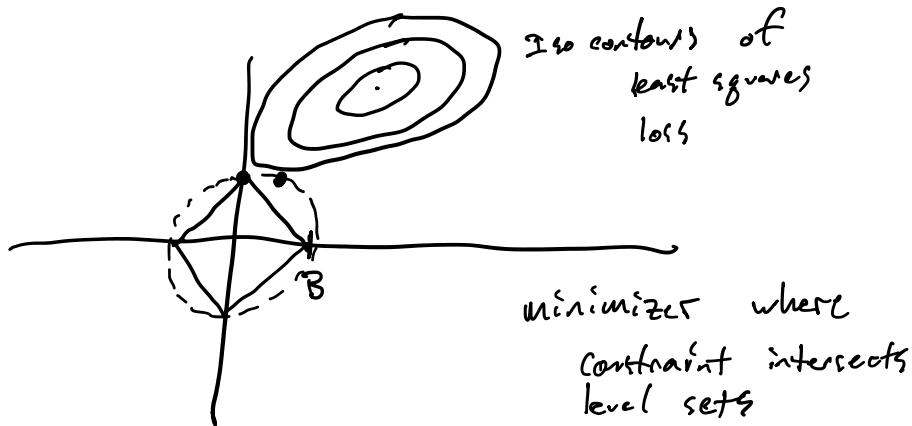
$$\min_w -\log p(y|w)$$

s.t. $\|w\|_1 \leq B$ \curvearrowleft parameter imposing tightness of penalty

Compare to ridge regression

$$\min_w -\log p(y|w)$$

s.t. $\|w\|_2 \leq B$



Issue Gradient undefined precisely where entries get sparse

Def Subgradient

$f: \mathbb{R}^n \rightarrow \mathbb{R}$, $g \in \mathbb{R}^n$ is a subgradient

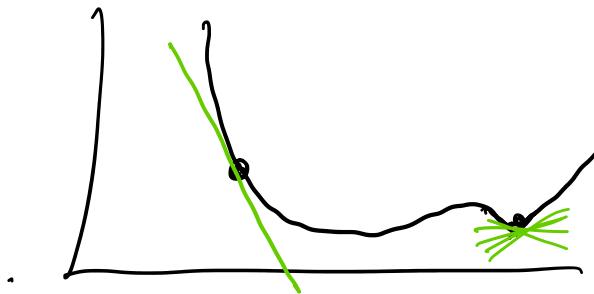
if for all $z \in \text{dom}(f)$

$$f(z) \geq f(x) + g^T(z-x)$$

\curvearrowleft think of like

Denote $g(x) = \partial f(x)$

first order Taylor series



At cont. pt, one
unique subdifferential
At a kink, many

Ex $f(x) = \text{abs}(x)$

$$\partial f(x) = \begin{cases} -1 & x < 0 \\ [-1, 1] & x = 0 \\ 1 & x > 0 \end{cases}$$

any number
between $[-1, 1]$

Back to lasso (Summary of results, see II.4.3 for details)

$$f(w) = \textcircled{1} \|Xw - y\|_2^2 + \lambda \textcircled{2} \|w\|,$$

One can show

$$\frac{\partial}{\partial w_d} \textcircled{1} = a_d w_d - c_d$$

$$a_d = \sum_{n=1}^N x_{nd}^2$$

interpret as residual $\rightarrow c_d = \sum_{n=1}^N x_{nd} (y_n - \vec{w}_{-d}^T \vec{x}_{n,-d})$

where $\vec{z}_{-d} = \langle z_1, \dots, z_{d-1}, z_{d+1}, \dots, z_m \rangle$

Taking subderivative to account for discontinuity

$$\partial_{w_d} (\textcircled{1} + \textcircled{2}) = (a_d w_d - c_d) + \lambda \partial_{w_d} \|w\|,$$

$$= \begin{cases} ad w_d - cd - 1 & \text{if } w_d < 0 \\ [-cd - \lambda, -cd + \lambda] & \text{if } w_d = 0 \\ ad w_d - cd + \lambda & \text{if } w_d > 0 \end{cases}$$

At $0 = 2w_d(0 + \lambda)$ we have 3 solutions

$$\hat{w}_d(c_d) = \begin{cases} (c_d + \lambda)/ad & c_d < -\lambda \\ 0 & c_d \in [-\lambda, \lambda] \\ (c_d - \lambda)/ad & c_d > \lambda \end{cases}$$

Compactly

$$\hat{w}_d = \text{Soft Threshold}\left(\frac{c_d}{ad}, \lambda/ad\right)$$

$$\text{SoftThreshold}(x, \delta) = \text{sign}(x) (|x| - \delta)_+$$

First introduction to Markov Chains

Consider time series data

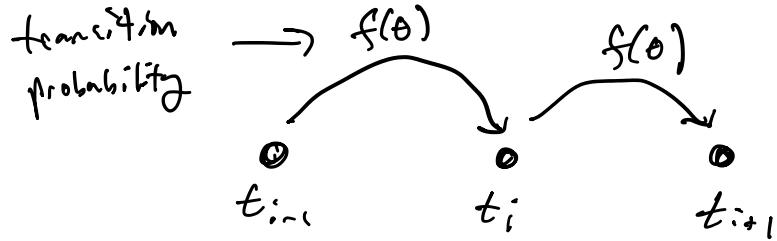
$$\vec{y}_d = (y_{t_1}^d, \dots, y_{t_N}^d)$$

By chain rule could write

$$\begin{aligned} P(\vec{y}) &= P(y_{t_1}) P(y_{t_2} | y_{t_1}) P(y_{t_3} | y_{t_2}, y_{t_1}) \dots \\ &= \prod_{i=1}^N P(y_{t_i} | y_{t_{i-1}}) \end{aligned}$$

For first-order Markov process, assume

$$P(y_{t_i} | y_{t_{i-1}}) = P(y_t | y_{t_{i-1}}) \\ = f(\theta)$$



So $P(\vec{y}_d) = P(y_1) \prod_{i=2}^n P(y_i | y_{i-1})$

↗
initial condition

Note as a special case a stochastic differential equation

$$dx = f(x) dt + g(x) dB_t$$

Euler-Maruyama $x^{n+1} - x^n = f(x^n) \Delta t + g(x^n) \xi$

$$\xi \sim N(0, \Delta t)$$

then $P(x^{n+1} | x^n) = N(x^n + \Delta t f(x^n), \Delta t g(x^n)^2)$