

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

Problem 2 Assume direct measurements of relevant DOFs z are unavailable

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

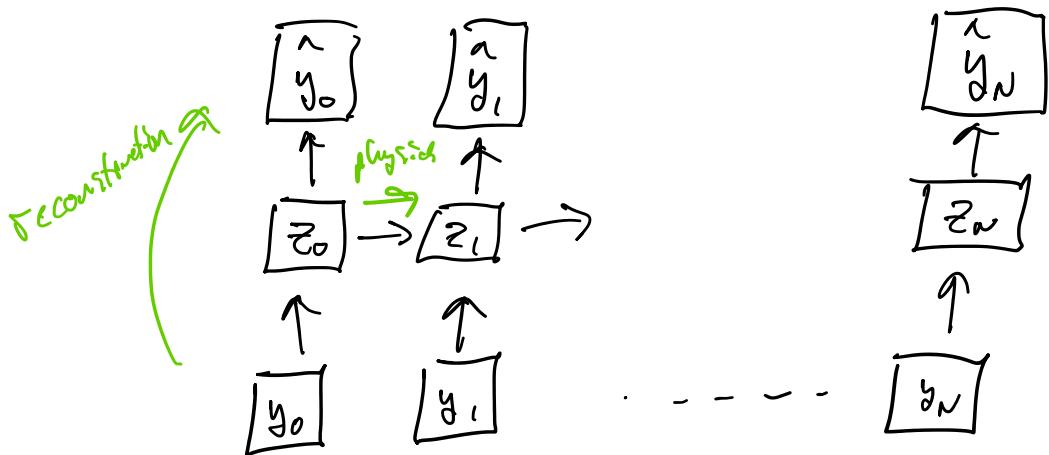
$$\begin{cases} \dot{z} = F(z, t | \theta_2) \\ z(t=0) = G(y(t=0)) \end{cases}$$

example Image-based simulation

Candidate architectures

$G \rightarrow$ auto encoders

$F \rightarrow$ Neural ODEs, universal differential
equations, 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 unknown coefficients

SPARSITY Want to only choose a few terms, so data reveals governing rule

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 $\ddot{\theta} = -k\theta$ is model

Can we infer θ coordinate from images?

Necessary Background

Revisiting overfitting (Chapters §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.
param ↑
complex

Typical Bayesian setting

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

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

def Maximum a posteriori (MAP) estimation

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

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

$$\hat{\theta} = \arg \max_{\theta} \left[\log p(\mathbf{x}|\theta) + \underbrace{\log p(\theta)}_{\substack{\text{likelihood} \\ \text{prior}}} \right]$$

note the unknown marginal likelihood $p(\mathbf{x})$
drops out

Bayesian regression revisited

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

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

$$\mathbf{M}_{MLE} = \mathbf{P}(\mathbf{x}) (\mathbf{P}(\mathbf{x}) \mathbf{P}(\mathbf{x})^T)^{-1} \mathbf{P}(\mathbf{x}) y$$

Ridge Regression

Adopt the prior

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

$$\mathcal{L} = \log p(y|\mathbf{x}, \mathbf{c}) p(\mathbf{c}) =$$

$$\log N(y; \mathbf{P}(\mathbf{x})^T \mathbf{c}, \sigma^2) N(\mathbf{c}; \mathbf{0}, \lambda^2 \mathbf{I})$$

$$= \sum_{i=1}^n -\frac{1}{2} \log \sigma^2 - \frac{1}{2} \left(\frac{\mathbf{P}(\mathbf{x})^T \mathbf{c} - y_i}{\sigma} \right)^2 - \frac{1}{2} \frac{\mathbf{c}^T \mathbf{c}}{\lambda^2}$$

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

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

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

\curvearrowleft 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_1 regression for sparsity

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

Goal

$$L_0 \text{ norm} \quad \|w\|_0 = \sum_{d=1}^D \mathbb{I}_{|w_d| > 0} \quad \text{minimization}$$

- Important for feature selection & dictionary

~ Important for feature selection & dictionary learning

Least Absolute shrinkage + selection operator
(Tibshirani 96)

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

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

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

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

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/

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

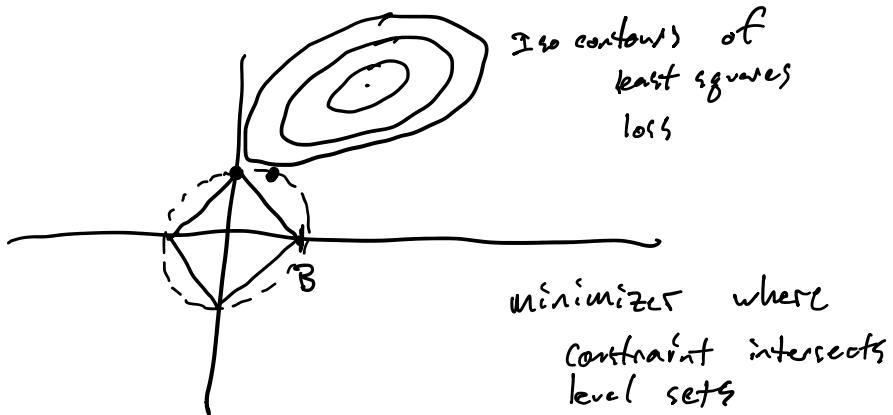
s.t. $\|w\|_1 \leq B$ \curvearrowleft parameter imposing tightness of penalty

Compare to ridge regression

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

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

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

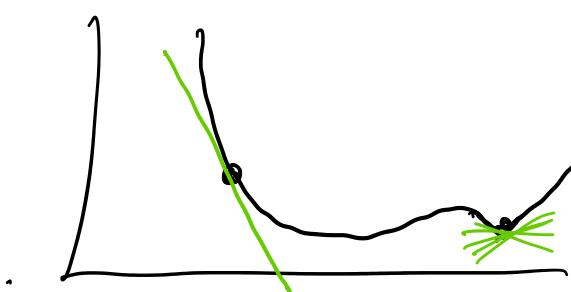
if for all $z \in \text{dom}(f)$

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

\curvearrowleft think of like

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

first order Taylor series



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

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

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

any number
between $[-1, 1]$

Back to Lasso (Summary of results, see 11.4.3 for details)

$$L(w) = \|Xw - y\|_2^2 + \lambda \|w\|,$$

One can show

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

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

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

$$\text{where } \hat{z}_{-d} = \langle z_1, \dots, z_{d-1}, z_{d+1}, \dots, z_N \rangle$$

Taking subderivative to acc't for discont.

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

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

At $0 = \partial_{w_d} (\textcircled{1} + \textcircled{2})$ we have 3 solutions

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

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

Compactly

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

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

First introduction to Markov Chains

Consider time series data

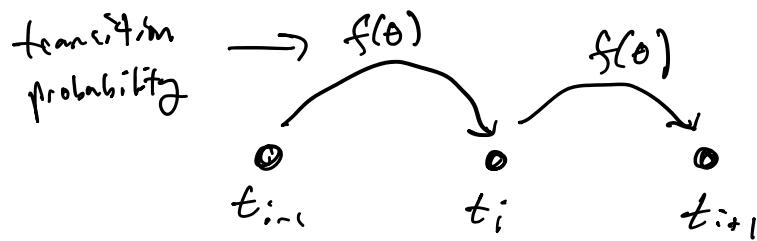
$$\vec{y}_d = (y_{t_1}^d, \dots, y_{t_N}^d)$$

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

$$\begin{aligned} P(y_{t_i} | y_{t_{i-1}}) &= P(y_i | y_{t_{i-1}}) \\ &= f(\theta) \end{aligned}$$



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

$$^{\text{70}} \quad p(\vec{y}_d) = p(y_1) \prod_{i=2}^n p(y_i | y_{i-1})$$

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

Exercise

Given data for harmonic oscillator

$$\theta(t) = \theta_0 + \sin \omega t$$

We aim to fit

Dynamics of the form

$$\ddot{z} = -\lambda z$$

\curvearrowleft single parameter to be fit

Discretize in time

$$\frac{z^{n+1} - 2z^n + z^{n-1}}{\Delta t^2} = -\lambda z^n$$

$$z^{n+1} = (2 + \Delta t^2 \lambda) z^n - z^{n-1}$$

Assume $z^n \sim N(M_n, \sigma_n^2)$ univariate gaussian

Recall identities

$$\text{If } A \sim N(M_A, \sigma_A^2)$$

$$B \sim N(M_B, \sigma_B^2)$$

$$\rightarrow \alpha A + \beta \sim N(\alpha M_A + \beta, \alpha^2 \sigma_A^2)$$

$$\rightarrow A + B \sim N(M_A + M_B, \sigma_A^2 + \sigma_B^2)$$

From
def of
exp. var
we didn't
show yet

$$\text{Then } p(z^{n+1} | z^n, z^{n-1}, \lambda) =$$

$$N\left(\underbrace{(2 + \Delta t^2 \lambda) M_n - M_{n-1}}_{M_{n,n-1}}, \underbrace{(2 + \Delta t^2 \lambda)^2 \sigma_n^2 + \sigma_{n-1}^2}_{C_{n,n-1}}\right)$$

Apply Markovian Assumption

$$p(z_0, \dots, z_n | \lambda) = p(z_0, z_1) \prod_{i=2}^n p(z_i | z_{i-1}, z_{i-2}, \lambda)$$

\downarrow indep of λ

$$L(\lambda) = -\log p(z | \lambda) = C + \sum_{i=1}^N \log p(z_i | z_{i-1}, z_{i-2}, \lambda)$$

Exercise 1

Given noisy data

$$z_n = \theta(t_n) + \varepsilon$$

$$\varepsilon \sim N(0, 0.01)$$

Infer λ

Exercise 2

Given indirect measurements

$$y_n = \begin{pmatrix} \cos z_n \\ \sin z_n \end{pmatrix}$$

Build Autoencoder

$$z_n \sim N(f_{\theta_1}(y_n), g_{\theta_2}(y_n))$$

Recall
ELBO

$$\mathcal{L}(\theta, \phi, \lambda | \mathcal{D}) = E_{q(z|x)} \left[\log \frac{p_{\theta}(x|z)}{q_{\phi}(z|x)} \right]$$

$$= E_{q(z|x)} \left[p(x|z) \right] - KL \left(q(z|x) || p(z) \right)$$

Choose prior:

$$p(\vec{z}) = p(z_0, z_1) \prod_{i=2}^n p(z_i | z_{i-1}, z_{i-2})$$

Assume separable posterior

$$q(\vec{z} | \vec{x}) = \prod_{i=0}^n q(z_i | x_i)$$

Then

$$KL(g(\vec{z}|\vec{x}) || p(z)) = \mathbb{E}_{g(\vec{z}|\vec{x})} \left[\log \frac{p(\vec{z})}{g(\vec{z}|\vec{x})} \right]$$

$$= \sum_{i=2}^N \mathbb{E}_{g(\vec{z}|\vec{x})} \left[\log p(z_i | z_{i-1}, z_{i-2}) \right] + C$$

$$- \sum_{i=0}^N \mathbb{E}_{g(\vec{z}|\vec{x})} \left[g(z_i | x_i) \right]$$

$$= \sum_{i=2}^N \mathbb{E}_{g(z_i|x_i)} \left[\log p(z_i | z_{i-1}, z_{i-2}) \right]$$

$$- \sum_{i=2}^N \mathbb{E}_{g(z_i|x_i)} \left[g(z_i | x_i) \right]$$

How can we compute these integrals?

Exercise Let $f(z) \sim N(\mu_1, \sigma_1^2)$
 $g(z) \sim N(\mu_2, \sigma_2^2)$

then

$$\int_{-\infty}^{\infty} \log(g(z)) dz = -\frac{1}{2} \sum \log 2\pi \sigma_{2,i}^2 + \frac{\sigma_{2,i}^2}{2}$$

$$\int_{R^N} f(z) \log(g(z)) dz = -\frac{1}{2} \sum_j \log 2\pi \sigma_{2,j}^2 + \frac{\sigma_{2,j}^2}{\sigma_{1,j}^2} + \frac{(M_{1,j} - M_{2,j})^2}{\sigma_{2,j}^2}$$