# ORIE 4742 - Info Theory and Bayesian ML

Chapter 8: Bayesian Regression

simplest "general" model for continuous data

- basis functions (for today - fixed basis)

- next class - infinite families of basis fns

March 29, 2021

Sid Banerjee, ORIE, Cornell

(eg- polynomial regression, but with no max degree)

'need some form of implicit regularization'

idea - Gaussian process

- Bayesian model selection

# normal-normal model for unknown $\mu$

- data $D = \{X_1, X_2, \ldots, X_n\} \in \mathbb{R}^n$
- model $\mathcal{M}$: $X_i$ i.i.d. from $\mathcal{N}(\mu, \tau)$, with unknown $\mu$, known $\tau = 1/\sigma^2$

*precision*

## normal-normal model

- likelihood: $p(D|\mu) \propto \exp\left(-\tau \sum_{i=1}^n (x_i - \mu)^2 / 2\right)$

- prior: $\mu \sim \mathcal{N}(m_\mu, 1/\tau_\mu) \propto \exp\left(-\tau_\mu(\mu - m_\mu)^2/2\right)$    $\tau_\mu, m_\mu = $ *prior hyperparameters*

- posterior: let $\overline{x} = \frac{1}{n}\sum_{i=1}^n x_i$ *(MLE)*, $m_D = \frac{n\tau \cdot \overline{x} + \tau_\mu \cdot m_\mu}{n\tau + \tau_\mu}$ *(shrinkage estimator)* and $\tau_D = n\tau + \tau_\mu$ *('precisions add')*

$$p(\mu|D) \sim \mathcal{N}(m_D, 1/\tau_D)$$

$$\mu = m_D + \frac{1}{\sqrt{\tau_D}} z_1, \qquad z_1, z_2 \sim N(0,1)$$

- posterior predictive distribution:

$$p(x|D) \sim \mathcal{N}(m_D, 1/\tau + 1/\tau_D)$$

$$X = \mu + \frac{1}{\tau} z_2$$

*aside*

if $X \sim \mathcal{N}(\mu, \sigma^2)$

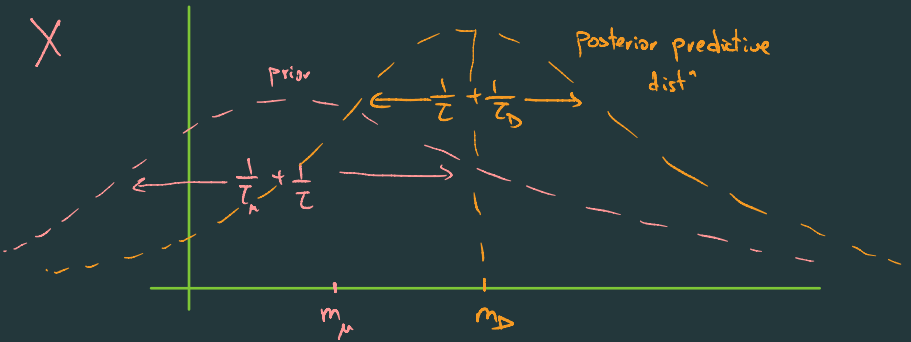$\Rightarrow X = \mu + \sigma z$

$z \sim N(0,1)$

$\mu$

posterior

prior

$\xleftarrow{\quad} \frac{1}{\tau_\mu + n\tau} \xrightarrow{\quad}$

$\xleftarrow{} \frac{1}{n\tau} \xrightarrow{}$

$\xleftarrow{\quad} \frac{1}{\tau_\mu} \xrightarrow{\quad}$

$m_\mu \qquad \tau_\mu \qquad m_D \qquad n\tau \quad \bar{x}$

$X$

prior

Posterior predictive dist$^n$

$\xleftarrow{} \frac{1}{\tau} + \frac{1}{\tau_D} \xrightarrow{}$

$\xleftarrow{\quad\quad} \frac{1}{\tau_\mu} + \frac{1}{\tau} \xrightarrow{\quad\quad}$

$m_\mu \qquad\qquad m_D$

# what is linear regression?

$\underline{\text{Data}}$ – $(x_1, t_1), (x_2, t_2), \ \text{---} \ , \ (x_n, t_n)$

observation  target

$\underline{\text{Model}}$ – $y(x) = \sum_{j=1}^{M} w_j \ \phi_j(x)$

weights  basis fns

quadratic regression

linear regression

$t_i$

$x_i$

$t(x) = y(x) + \varepsilon$  iid noise $\varepsilon \sim N(0, 1/\beta)$

assume $\phi_1(x) = 1$  noise precision

$\underline{E_g}$ – linear regression – $t(x) = W_1 + W_2 x + \varepsilon$

$(\text{degree } 3)$ polynomial regression – $t(x) = W_1 + W_2 x + W_3 x^2 + W_4 x^3 + \varepsilon$

$\phi \equiv (1, x, x^2, x^3)$

# basis functions (from Bishop Ch 6)



Polynomial basis

$1, x, x^2, x^3, x^4, \ldots$

$\phi_j(x) = x^{j-1}$ (Taylor series)

- $\phi_j(x) = \sin(\omega_j x + \mu_j)$ — Fourier basis (Fourier series)

- Wavelet basis

Gaussian basis fn

$$\phi_j(x) = e^{-(x-\mu_j)^2/s_j}$$

location
scale

Sigmoidal basis fn

$$\phi_j(x) = \frac{1}{1 + e^{-(x-\mu_j)/s_j}}$$

# regression: the frequentist view

$$t_i = \sum_{j=1}^{M} w_i \phi_i(x_i) + \varepsilon_i \quad , \quad \varepsilon_i \sim N(0, 1/\beta) \text{ , iid}$$

Design matrix 
$$\phi = \begin{pmatrix} \phi_1(x_1) & \phi_2(x_1) & \cdots & \phi_M(x_1) \\ \phi_1(x_2) & \phi_2(x_2) & \cdots & \phi_M(x_2) \\ \vdots & & & \\ \phi_1(x_n) & \phi_2(x_n) & \cdots & \phi_M(x_n) \end{pmatrix}$$

$n \times M$ matrix

$$\phi(x) = (\phi_1(x), \phi_2(x), \ldots, \phi_M(x))^T$$

$$w = (w_1, w_2, \ldots, w_M)^T$$

$$t = (t_1, t_2, \ldots, t_n)^T$$

$(\phi, t)$ are a sufficient statistic for the data $\left(\begin{smallmatrix}\text{under this}\\\text{model}\end{smallmatrix}\right)$

likelihood 
$$p(D | w, M) \propto \exp\left(-\sum_{i=1}^{n} \frac{\hat{\beta}}{2}(t_i - w^T \phi(x))^2\right)$$
model

$(\text{assuming } \beta \text{ is known})$

# Frequentis regression

- maximize $\mathcal{L}_D(w)$ $\Longleftrightarrow$ maximize $l(w) = \log \mathcal{L}(w)$

$\Longleftrightarrow$ minimize $\displaystyle\sum_{i=1}^{n} \left( t_i - w^T \phi(x_i) \right)^2$

ie - Least squares!

- $W_{MLE} = \underbrace{\phi^{\dagger}}_{\substack{\text{Moore-Penrose} \\ \text{pseudoinverse}}} t = \underbrace{\left( \phi^T \phi \right)^{-1} \phi^T t}_{\substack{\text{MLE estimator for regression} \\ \text{coeffs}}}$
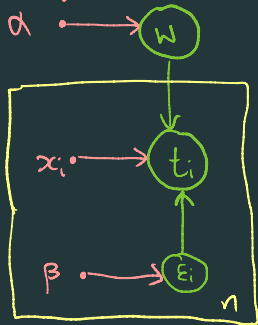
# Bayesian linear regression

**Model** -
$$t_i = \sum_{j=1}^{M} W_j \, \phi_j(x_i) + \varepsilon_i$$

- Now $W_1, W_2, \ldots, W_M$ are random (but common to all data)

$$\varepsilon_i \sim N(0, 1/\beta) \,, \; iid$$

prior hyperparameters



$$\bullet \quad W = \begin{pmatrix} W_1 \\ W_2 \\ \vdots \\ W_M \end{pmatrix} \sim N\left(0, T_0^{-1}\right)$$

↑
precision matrix

$$\left( i.e. \; T_0 = \Sigma^{-2} \right)$$

$$\underline{eg} - T_0 = \alpha^{-1} I$$

$$\Rightarrow \sigma^2 = 1/\alpha \,, \text{ and } W_i \sim N(0, \sigma^2), \; iid$$

# Bayesian linear regression *(generalizes the normal-normal model to M dimensions)*

- data $D = \{(t_1, x_1), (t_2, x_2), \ldots, (t_N, X_N)\} \in \mathbb{R}^n$
- model $\mathcal{M}$: $t_i = \sum_{j=0}^{M-1} W_j \phi(x_i) + \epsilon_i$, where $\epsilon_i \sim \mathcal{N}(0, \beta^{-1})$ ← *known precision*

**Bayesian linear regression model** - *hyper params* - $\beta$, $\alpha$, $(M, (scale, loc) \text{ for bases fns})$

- likelihood: $p(D|W) \propto \exp\left(-\beta \sum_{i=1}^{N}(x_i - W^\intercal \phi(x_i))^2 / 2\right)$

- prior: $W \sim \mathcal{N}(0, \alpha^{-1}I)$ ← $m_w$, *ie,* $W_j \sim \mathcal{N}(0, 1/\alpha)$

- posterior: let $m_D = T_D^{-1} \beta \Phi^\intercal t$ and $T_D = \beta \Phi^\intercal \Phi + \alpha I$
  
  $\underbrace{T_D^{-1} \beta \Phi^\intercal t}_{\text{pseudo inverse}}$   $\underbrace{\beta \Phi^\intercal \Phi}_{\text{data precision}} + \underbrace{\alpha I}_{\text{prior precision}}$

  $$p(W|D) \sim \mathcal{N}\left(m_D, T_D^{-1}\right)$$

Recall - $W_{MLE} = \left(\Phi^\intercal \Phi\right)^{-1} \Phi^\intercal t$

$m_D = \left(\Phi^\intercal \Phi + \dfrac{\alpha}{\beta} I\right)^{-1} \Phi^\intercal t$
$\underbrace{\qquad\qquad}_{\text{regularizer}}$

Note - even though $W_i$ were indep in prior, they are dependent conditioned on data (explaining away...)

# Bayesian linear regression: example (Bishop chapter 3)



Model — $t_i = w_0 + w_1 x_i + \varepsilon_i$

$$W = \begin{pmatrix} w_0 \\ w_1 \end{pmatrix} \sim \mathcal{N}(0, \alpha^{-1} I), \quad \varepsilon_i \sim \mathcal{N}(0, 1/\beta)$$

$$\Phi = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ & \vdots \\ 1 & x_n \end{pmatrix}, \quad t = \begin{pmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{pmatrix}$$

- $m_D = \left( \Phi^T \Phi + \dfrac{\alpha}{\beta} I \right)^{-1} \Phi^T t$

$$T_D = \beta \left( \Phi^T \Phi + \dfrac{\alpha}{\beta} I \right)$$

- As $n \nearrow \alpha$, $T_D \searrow 0$

$$m_D \to \begin{pmatrix} -0.3 \\ 0.1 \end{pmatrix}$$

(In figure, labels: likelihood, prior/posterior, data space; $w_0, w_1$ are correlated; true value)

ground truth: $f(x) = 0.1x - 0.3$

## Bayesian linear regression

- data $D = \{(t_1, x_1), (t_2, x_2), \ldots, (t_N, X_N)\} \in \mathbb{R}^n$
- model $\mathcal{M}$: $t_i = \sum_{j=0}^{M-1} W_j \phi(x_i) + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \beta^{-1})$

### Bayesian linear regression model

- likelihood: $p(D|W) \propto \exp\left(-\beta \sum_{i=1}^{N}(x_i - W^\intercal \phi(x_i))^2/2\right)$

- prior: $W \sim \mathcal{N}(0, \alpha^{-1}I)$

- posterior: let $m_D = T_D^{-1}\beta\Phi^\intercal t$ and $T_D = \beta\Phi^\intercal\Phi + \alpha I$

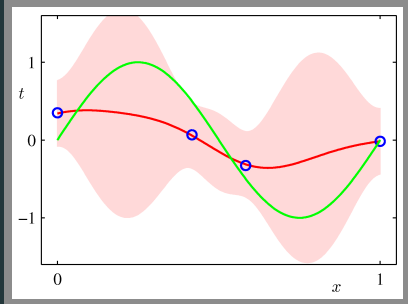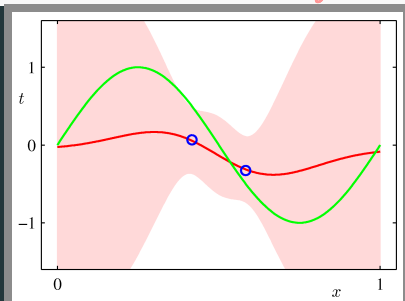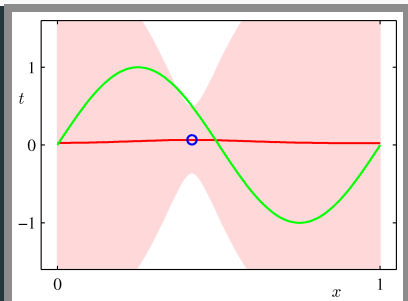$$p(W|D) \sim \mathcal{N}\left(m_D, T_D^{-1}\right)$$

- posterior predictive distribution: *ie., what is $p(t|D)$ for new $x$*

$$p(t|D) \sim \mathcal{N}\left(m_D^\intercal \phi(x), \underbrace{\beta^{-1} + \phi(x)^\intercal T_D^{-1}\phi(x)}\right)$$

*variances add up, depends on $x$*

# Bayesian linear regression: posterior prediction

# Bayesian linear regression: posterior sampling