<u>Last 4 classes</u>
- Probabilistic graphical models
- MC MC / Monte Carlo

# ORIE 4742 - Info Theory and Bayesian ML

Bayesian Regression $\left(\begin{array}{l}\text{today - fixed basis functions} \\ \text{next class - Gaussian processes}\end{array}\right)$
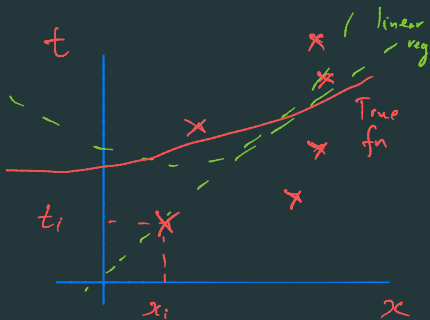
April 23, 2020

Sid Banerjee, ORIE, Cornell

# what is linear regression?

Data - $(x_1, t_1), (x_2, t_2), \ldots, (x_N, t_N)$

↑ observations ↑ target



quadratic reg

linear reg

True fn

- ## Model

- $y(x) = \sum_{j=0}^{M-1} w_j \, \phi_j(x)$
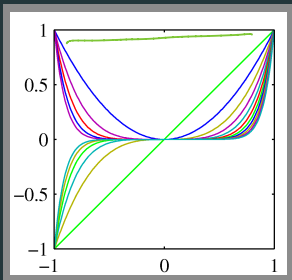
↑ regression coefficient

↑ basis vectors

- $t(x) = y(x) + \varepsilon, \quad \varepsilon \sim N(0, 1/\beta) - \text{Noise}$
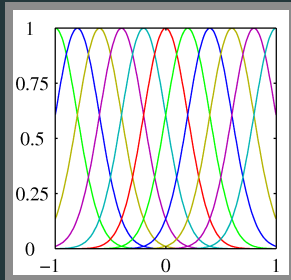
← noise precision

frequentist view of regression

- Assume $\phi_0(x) = 1 \quad (w_0 \equiv \text{constant, 'bias'})$
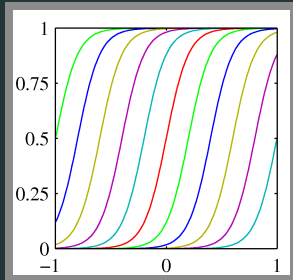
# basis functions



Polynomial basis fns

$$\phi_j(x) = x^j$$

Gaussian basis fn

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

location parameter    scale parameter

Sigmoidal basis fn

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

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

- Wavelet basis

# regression: the frequentist view $y(x) = w_0 + w_1 x$

$(M)$ $\quad t(x) = \sum_{j=0}^{M-1} w_j \phi_j(x) + \varepsilon, \quad \varepsilon \sim N(0, 1/\beta)$

- **design matrix**

$$\Phi = \begin{pmatrix} \phi_0(x_1) & \phi_1(x_1) & \cdots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \cdots & \phi_{M-1}(x_2) \\ \vdots & & & \\ \phi_0(x_N) & \phi_1(x_N) & \cdots & \phi_{M-1}(x_N) \end{pmatrix}$$

$N \times M$ matrix

$$D = \overline{\Phi}, \quad t = (t_1, t_2, \ldots, t_N)^T$$

$N \times 1$ vector

$M \times 1$ vector $\quad w = (w_0, w_1, \ldots, w_{M-1})$

$\phi(x_i) = (\phi_0(x) \ldots \phi_{M-1}(x))$

$\underbrace{\qquad\qquad\qquad\qquad\qquad}$

Sufficient statistic of the data

- **likelihood** $\quad P(D|M) \propto \exp\left(-\sum_{i=1}^{N} \beta \frac{(t_i - w^T \phi(x))^2}{2}\right)$

- **maximum likelihood** - $w_{ML} = \underset{\text{dagger} \rightarrow}{\Phi^\dagger} t, \quad \Phi^\dagger = (\Phi^T \Phi)^{-1} \Phi^T$

pseudo inverse $\qquad M \times N$ matrix

Eg - linear regression (frequentist)

$$t = w_0 + w_1 x + \varepsilon$$

observed data ← (above $t$)

unknown params (below $w_0, w_1$)

$\left( t = y(x) + \varepsilon, \; y(x) = w_0 + w_1 x \right)$

noise $N(0, 1/\beta)$

$$\Phi = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \\ 1 & x_N \end{pmatrix}, \quad t = \begin{pmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{pmatrix}, \quad w = \begin{pmatrix} w_0 \\ w_1 \end{pmatrix}$$

- $$w^{ML} = \underbrace{(\Phi^T \Phi)^{-1} \Phi^T}_{A(x_1, x_2, \ldots, x_N)} \, t$$

Alternate
- choose $w_0, w_1$ to
  minimize $\sum_{i=1}^{N} (t_i - w_0 - w_1 x_i)^2$
- ie, LS estimate

- output - $\boxed{y(x) = w_0^{ML} + w_1^{ML} x}$

**Bayesian linear regression**

$$P(t \mid w) \sim N\left(w^T \phi(x), \tfrac{1}{\beta}\right)$$

<u>Model</u> -

$$t_i = \sum_{j=0}^{M-1} w_j \, \phi_j(x_i) + \varepsilon_i$$

'unknown' $\equiv$ random variables

— $\varepsilon_i \sim N\left(0, \tfrac{1}{\beta}\right) \equiv$ iid for each $(x_i, t_i)$

— $w = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_{M-1} \end{pmatrix} \sim N\left(0, T_0^{-1}\right)$  (prior)

$\begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$  'precision' matrix

$Eg - \alpha^{-1} I$

i.e, $w_j \sim N(0, \tfrac{1}{\alpha})$, iid $\forall j$

— $\alpha, \beta \equiv$ model hyper parameters (fixed)

## 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 $\tau = 1/\sigma^2$
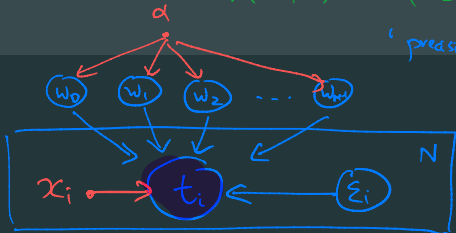
### normal-normal model

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

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

- posterior: let $\overline{x} = \frac{1}{n}\sum_{i=1}^{n} x_i$, $m_D = \frac{n\tau \cdot \overline{x} + \tau_\mu \cdot m_\mu}{n\tau + \tau_\mu}$ and $\tau_D = n\tau + \tau_\mu$

  $\underbrace{\phantom{xxxxxxxx}}$
  
  *empirical mean*
  
  (ML estim for $\mu$)

  $p(\mu|D) \sim \mathcal{N}(m_D, 1/\tau_D)$ $\quad$ 'shrinkage' estimator - $\tau \overline{x} + (1-\tau)m_\mu$

- posterior predictive distribution:

  noise added to X by model

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

  'noise' in parameter $\mu$

# Bayesian linear regression

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

## Bayesian linear regression model

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

- prior: $W \sim \mathcal{N}(0, \alpha^{-1}I)$ $\left(i.e.,\ W_j \sim \mathcal{N}\left(0, \frac{1}{\alpha}\right),\ iid\right)$

$$m_D = T_D^{-1}\beta \sum_{i=1}^{N} \phi(x_i) t_i$$
$$\underset{M \times M}{}\quad \underset{M \times 1}{}$$

- posterior: let $m_D = \underset{M\times 1}{T_D^{-1}}\underset{M_D \times M}{\beta}\underbrace{\Phi^\top t}_{M \times 1}$ and $T_D = \underset{M \times M}{\beta\Phi^\top\Phi + \alpha I}$

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

'$precisn$' inverse covariance matrix



$\alpha$

$W_0$  $W_1$  $W_2$  $\cdots$  $W_M$

$x_i \longrightarrow t_i \longleftarrow \xi_i$  $N$

- Note
$\{W_i\}_{i=0}^{M-1}$ initially indep,
but dependent given $t$.

# Bayesian linear regression: example   (from Bishop Ch 3)



| likelihood | prior/posterior | data space |

distribution over fns →

**model** - $t_i = W_0 + W_1 x_i + \varepsilon_i$

$$\begin{pmatrix} W_0 \\ W_1 \end{pmatrix} \sim N(0, \alpha^{-1} I), \quad \varepsilon_i \sim N(0, 1/\beta)$$

• $y(x) = -0.3 + 0.1 x, \quad t_i = y(x)_i + \varepsilon_i$

---

* As $N$ increases

$$T_D \searrow 0$$

$$M_D \rightarrow \text{true} \begin{pmatrix} W_0 \\ W_1 \end{pmatrix}$$

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 $\mathcal{M}$: $t_i = \sum_{j=0}^{M-1} W_j \phi(x_i) + \epsilon_i$, where $\epsilon_i \sim \mathcal{N}(0, \beta^{-1})$

### 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: $\left(i.e, \ p(t|x,D)\ \right)$

$$p(t|D) \sim \mathcal{N}\left(m_D^\intercal\phi(x), \beta^{-1} + \phi(x)^\intercal T_D^{-1}\phi(x)\right) \xleftarrow{\text{Variance as}} \text{fn of } x$$

$i.e - W = m_D + Z, \ Z \sim \mathcal{N}(0, T_D), \ t = W^\intercal\phi(x) + \varepsilon_i$

$\Rightarrow t = m_D^\intercal\phi(x) + \underline{Z^\intercal\phi(x) + \varepsilon_i} \sim \mathcal{N}(0, \frac{1}{\beta} + \phi(x)^\intercal T_D^{-1}\phi(x))$
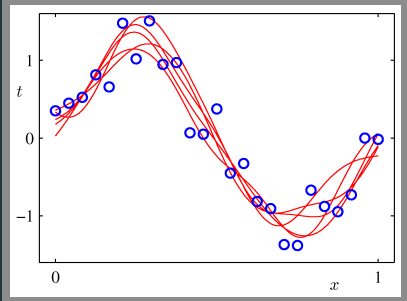
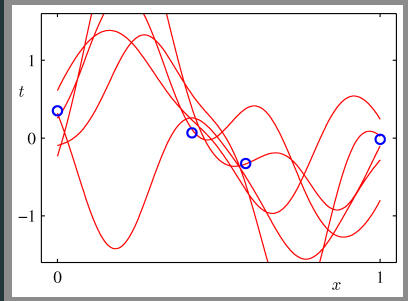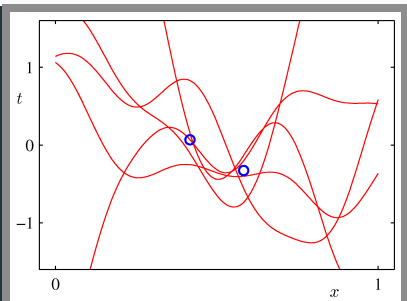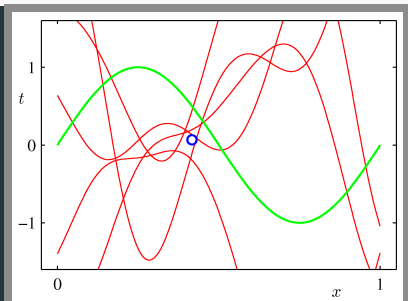# Bayesian linear regression: posterior prediction

Bishop Ch 3
- ground truth- $\sin 2\pi x$
- basis fn- Gaussian, M=10



$t(x_i) = y(x) + \varepsilon_i$

$y(x) = \sin(2\pi x)$

model

$t(x) = \sum_{j=0}^{9} \phi(x_i) w_j + \varepsilon_i$

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

# Bayesian linear regression: posterior sampling

**the 'equivalent' kernel** $\quad \left( \text{distance fn defined by data} \right)$

• given $D = \{(t_i, x_i)\}$, posterior mean $= t(x) = \sum_{i=1}^{N} t_i \, k(x, x_i)$

- 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 $\epsilon_i \sim \mathcal{N}(0, \beta^{-1})$
- prior: $W \sim \mathcal{N}(0, \alpha^{-1} I)$
- posterior: let $m_D = T_D^{-1} \beta \Phi^\mathsf{T} t$ and $T_D = \beta \Phi^\mathsf{T} \Phi + \alpha I$, then

$$t(x|D) = m_D^\mathsf{T} \phi(x) + \epsilon_D$$

↙ noise in model

where $\epsilon_D \sim \mathcal{N}(0, \beta^{-1} + \Phi^\mathsf{T} T_D^{-1} \Phi^\mathsf{T})$

← noise in params $W$ $\qquad T_D^{-1}$

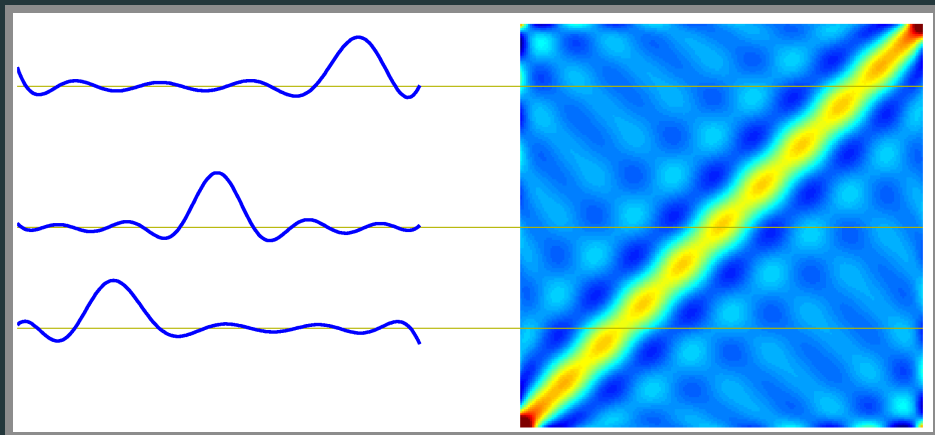alternately, $y(x|D) = \sum_{n=1}^{N} k(x, x_n) t_n$, where $k(x, y) = \beta \phi(x)^T S_D \phi(y)$

$$y(x \mid D) = m_D^\mathsf{T} \phi(x) = \beta \left( T_D^{-1} \Phi^\mathsf{T} t \right)^\mathsf{T} \begin{pmatrix} \phi_0(x) \\ \phi_1(x) \\ \vdots \\ \phi_{M-1}(x) \end{pmatrix}$$

$$= t^\mathsf{T} \left( \beta \Phi \, T_D^{-1}{}^\mathsf{T} \phi(x) \right)$$

$\underline{(AB)^\mathsf{T} = B^\mathsf{T} A^\mathsf{T}}$ $\qquad\qquad\qquad \underset{(\phi(x_i))_{i,j}}{}$

Polynomial kernel

Sigmoidal kernel

# The kernel view



$k(z, z_4)$

$k(z, z_5)$

$k(z, z_3)$

$k(z, z_2)$