

Last week (Generative models for discrete)

- The Dirichlet model (for multiclass data)
- The Naive Bayes classifier

Today

- Generative models for continuous data
- Gaussian - Gaussian, Gaussian - Gamma models
- Bayesian regression
- Model selection & the Bayesian Occam's razor
- Gaussian Processes

generative models for continuous data

Example - Regression (Bayesian)

Data: (x_i, t_i) , $i = 1, 2, \dots, n$

Model: $t = \underbrace{w_0 + w_1 x + w_2 x^2}_{\substack{\text{polynomial} \\ \text{regression model}}} + \underbrace{\varepsilon}_{\substack{\text{noise, } N(0, \sigma^2)}}$



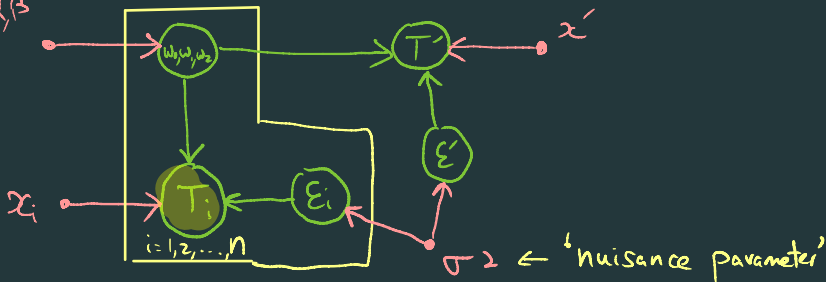
Q: find the 'best' degree 2 polynomial that 'approximates' data

Idea: Assume w_0, w_1, w_2 are random variables, learn from data via Bayesian update

Bayes Net for regression (for details, see Bishop Ch 8)

hyperparams for prior

α, β



$$\bullet T' \perp\!\!\!\perp T_i \mid w$$

- $\bullet T'$ not independent of data $\{T_i\}$ if not conditioned on w

continuous data and Gaussian priors

- data $D = \{X_1, X_2, \dots, X_n\} \in \mathbb{R}^n$, $X_i \in \mathbb{R}$ (1-dim data)
- model \mathcal{M} : X_i generated i.i.d. from $\mathcal{N}(\mu, \sigma^2)$ distribution

Gaussian prior

- $x \in \mathbb{R}$, parameters: $\Theta = (\mu, \sigma)$

- pdf: $\mathbb{P}[X_1 = x_1, X_2 = x_2, \dots, X_n = x_n] \propto \frac{1}{(\sigma^2)^{n/2}} \exp\left(\frac{-\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}\right)$

- normalizing constant: $(2\pi)^{-n/2}$

Conditioned on μ, σ^2
Problem: functional form of μ
and σ^2 are different
 $\mu \in \mathbb{R}, \sigma^2 \in \mathbb{R}_+$

3 options:

1. μ unknown, σ^2 known most important
2. σ^2 unknown, μ known
3. μ unknown, σ^2 unknown

notation: define precision $\tau = \frac{1}{\sigma^2}$

case 1: unknown μ

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

normal-normal model

μ
↑
random variable

Gaussian likelihood fn

- likelihood:

$$\mathcal{L}(\mu) = p(D|\mu) \propto \tau^{n/2} \exp\left(-\tau \sum_{i=1}^n (x_i - \mu)^2 / 2\right) = c_1 e^{-c_2(\mu - c_3)^2}$$

- prior parameter: $\Theta_0 = (m_\mu, 1/\tau_\mu)$ (mean, precision for μ)

- **Gaussian prior** for μ : $\mu \sim \mathcal{N}(m_\mu, \tau_\mu)$, where $\tau_\mu = 1/\text{Var}(\mu)$

hyperparameters of prior

$$p(\mu|m_\mu, \tau_\mu) \propto \tau_\mu^{1/2} \exp\left(-\tau_\mu(\mu - m_\mu)^2 / 2\right)$$



normal-normal model: posterior

$$P(\mu | D) \propto \underbrace{P(\mu)}_{\text{prior}} \cdot \underbrace{\mathcal{L}_D(\mu)}_{\text{likelihood}}$$

$$\begin{aligned} \cdot A &= \tau_{\mu} m_{\mu} + \tau \sum_{i=1}^n x_i \\ \cdot \bar{x} &= \frac{1}{n} \sum_{i=1}^n x_i = \begin{matrix} \text{sample} \\ \text{mean} \\ \text{(MLE)} \end{matrix} \end{aligned}$$

$$\cdot m_D = \frac{\tau_{\mu} m_{\mu} + n\tau \bar{x}}{\tau_{\mu} + n\tau}$$

$$\cdot \tau_D = \tau_{\mu} + n\tau$$

$$\propto \exp\left(\underbrace{-\frac{\tau_{\mu} (\mu - m_{\mu})^2}{2}}_{\text{want } e^{-c(\mu-c)^2}}\right) \exp\left(-\frac{\tau}{2} \sum_{i=1}^n (x_i - \mu)^2\right)$$

$$\propto \exp\left(\frac{-(\tau_{\mu} + n\tau)\mu^2 + 2A\mu + B}{2}\right)$$

$$\propto \exp\left(-\frac{\tau_D}{2} (\mu - m_D)^2\right)$$

normal-normal model: posterior

$$\Rightarrow \text{posterior} \equiv \mu \sim \mathcal{N}(m_D, \tau_D)$$

• posterior mean

$$m_D = \frac{\tau_\mu m_\mu + n\tau \bar{x}}{\tau_\mu + n\tau}$$
$$= w_{\text{prior}} m_\mu + (1 - w_{\text{prior}}) \underbrace{\bar{x}}_{\substack{\text{MLE} \\ \text{estimate}}}$$

'shrinkage estimator'

• posterior precision

$$\tau_D = \tau_\mu + n\tau$$

'precision on mean adds up under conditioning'

normal-normal model: posterior predictive distribution

• Prior on $\mu \sim \mathcal{N}(m_\mu, \tau_\mu)$, posterior given data $\mu_D \sim \mathcal{N}(m_D, \tau_D)$

$$\begin{aligned} \cdot X_{|\mu} &\sim \mathcal{N}(\mu, \tau) = \underbrace{\mu}_{\sim \mathcal{N}(m_D, \tau_D)} + \frac{1}{\tau} \mathcal{N}(0, 1) \\ &= m_D + \frac{1}{\tau_D} \mathcal{N}(0, 1) + \frac{1}{\tau} \mathcal{N}(0, 1) \\ &\sim \mathcal{N}\left(m_D, \underbrace{\frac{1}{\tau_D} + \frac{1}{\tau}}_{\sigma_D^2 + \sigma^2}\right) \end{aligned}$$

'posterior over $X \sim$ Gaussian, mean \equiv convex comb of m_μ, \bar{x}
variance \equiv sum of σ_D^2, σ^2

ie - 'uncertainties add up for posterior prediction'

normal-normal model: posterior predictive distribution

- data $D = \{X_1, X_2, \dots, X_n\} \in \mathbb{R}^n$
- model \mathcal{M} : X_i i.i.d. from $\mathcal{N}(\mu, \tau)$, with unknown μ , known $\tau = 1/\sigma^2$
- thus we have $Z_1, Z_2 \sim \mathcal{N}(0, 1)$ iid

$$X_i = \mu + \sigma Z_1$$

$$\mu = m_\mu + \sigma_\mu Z_2 = \mu + \sigma Z_1 + \sigma_\mu Z_2$$

$$\Rightarrow E[X_i] = E[\mu] = m_D, \quad \text{Var}(X_i) = \sigma^2 + \sigma_D^2$$

normal-normal model for unknown μ (see Bishop Ch 2, sec 3?)

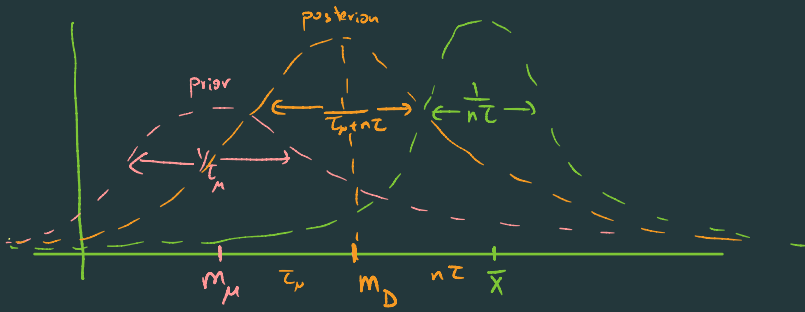
- data $D = \{X_1, X_2, \dots, X_n\} \in \mathbb{R}^n \rightarrow$ suff stat $x_{MLE} = \bar{x} = \sum x_i / n$
- model \mathcal{M} : X_i i.i.d. from $\mathcal{N}(\mu, \tau)$, with unknown μ , known $\tau = 1/\sigma^2$

normal-normal model

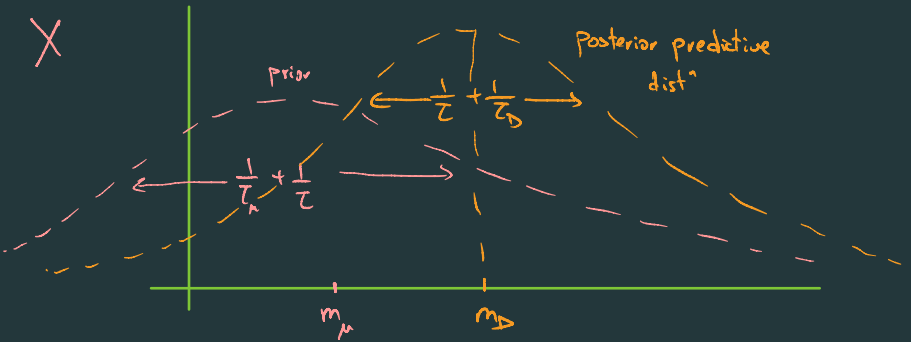
- likelihood: $p(D|\mu) \propto \exp(-\tau \sum_{i=1}^n (x_i - \mu)^2 / 2)$
- prior: $\mu \sim \mathcal{N}(m_\mu, 1/\tau_\mu) \propto \exp(-\tau_\mu (\mu - m_\mu)^2 / 2)$
- posterior: let $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$, $m_D = \frac{n\tau \cdot \bar{x} + \tau_\mu \cdot m_\mu}{n\tau + \tau_\mu}$ and $\tau_D = \frac{n\tau + \tau_\mu}{\text{precisions add for inference}}$
$$p(\mu|D) \sim \mathcal{N}(m_D, 1/\tau_D)$$
- posterior predictive distribution:

$$p(x|D) \sim \mathcal{N}(m_D, \underbrace{1/\tau + 1/\tau_D}_{\text{variances add for prediction}})$$

μ



X

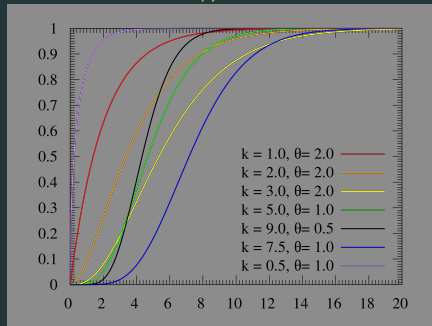
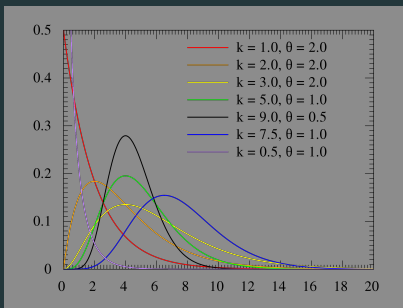


the gamma distribution

gamma distribution

- $x \in (0, \infty)$, parameters: $\Theta = (\alpha, \beta) \in \mathbb{R}^+$ ('shape, rate')
- pdf of $\text{Gamma}(\alpha, \beta)$: $p(x) \propto x^{\alpha-1} e^{-\beta x}$
- normalizing constant: $\frac{1}{Z(\alpha, \beta)} = \frac{\beta^\alpha}{\Gamma(\alpha)}$

Eg - If $X_i \sim \text{Exp}(\lambda)$ iid, then $X_1 + X_2 \sim \text{Gamma}$ with mean $2/\lambda$



case 2: unknown σ

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

normal-gamma model

- **likelihood**:

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

- prior parameters: $\Theta_0 = (\alpha, \beta)$
- **gamma prior** for τ : $\tau \sim \text{Gamma}(\alpha, \beta)$

$$p(\tau|\alpha, \beta) \propto \tau^{\alpha-1} e^{-\beta\tau}$$

normal-gamma model: posterior

$$P(\tau | D) \propto \underbrace{\tau^{\alpha-1} e^{-\beta\tau}}_{\text{prior}} \underbrace{\tau^{n/2} e^{-\frac{\tau}{2} \sum (x_i - \mu)^2}}_{\text{Likelihood}}$$

$$\propto \tau^{\alpha + \frac{n}{2} - 1} e^{-\tau \left(\beta + \frac{\sum (x_i - \mu)^2}{2} \right)}$$

$$= \text{Gamma} \left(\alpha + \frac{n}{2}, \beta + \frac{\sum (x_i - \mu)^2}{2} \right)$$

normal-gamma model: posterior predictive distribution

$$p(x | \alpha, \beta, D) = \int_0^{\infty} \underbrace{p(\tau | \alpha, \beta, D)}_{\text{gamma}} \cdot \underbrace{p(x | \tau)}_{\text{Gaussian}} d\tau$$

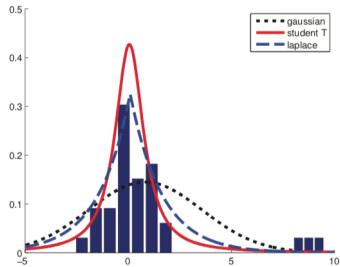
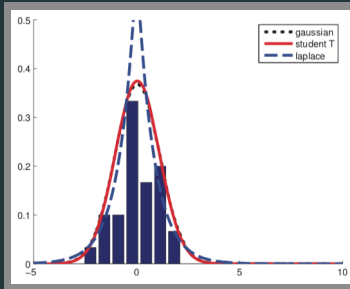
$$= \frac{1}{Z} \underbrace{\frac{1}{\left(1 + \frac{(x - \mu)^2}{2\beta}\right)^{d_D + 1/2}}}_{\text{student's t-dist}}$$

the Student-t distribution

(‘naturally robust’ distribution)

Student-t distribution

- $x \in \mathbb{R}$, parameter: $\mu \in \mathbb{R}, \nu > 0$ (mean, ‘degrees of freedom’)
- pdf of student-t(μ, ν): $p(x) \propto \left(1 + \frac{(x-\mu)^2}{\nu}\right)^{\frac{\nu+1}{2}}$
- normalizing constant: $\frac{1}{Z(\mu, \nu)} = \frac{\Gamma(\nu+1)/2}{\sqrt{\nu\pi}\Gamma(\nu/2)}$



robustness of student-t to outliers

normal-gamma model for unknown τ

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

normal-gamma model

- **likelihood**: $p(D|\theta) \propto \exp(-\tau \sum_{i=1}^n (x_i - \mu)^2 / 2)$
- **prior** for τ : $\tau \sim \text{gamma}(\alpha, \beta)$
- **posterior**: let $\alpha_D = \alpha + \frac{n}{2}$ and $\beta_D = \beta + \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2$

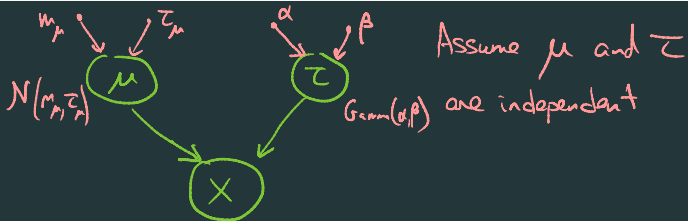
$$p(\tau|D) \sim \text{gamma}(\alpha_D, \beta_D)$$

- **posterior predictive distribution**:

$$p(x|D) \sim \text{student-t}$$

case 3: unknown μ and σ^2

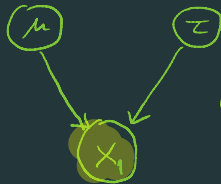
Idea 1



Now given data X_1 , what happens to μ, τ ?

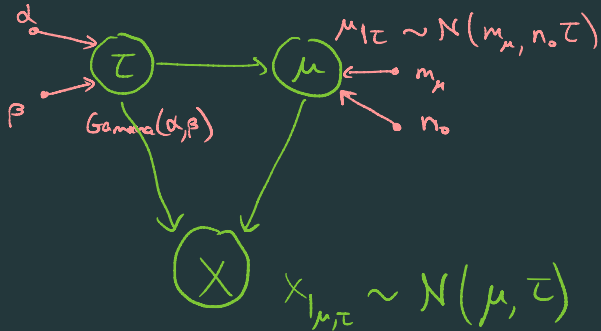
want posterior to be

this is not true!



Conditioned on X_1 , μ and τ are no longer indep ('explaining away')

case 3: unknown μ and σ^2



Bayesian update for this is known in closed form ('correct' conjugate prior)