

Lecture Slides for INTRODUCTION TO MACHINE LEARNING 3RD EDITION

ETHEM ALPAYDIN © The MIT Press, 2014

alpaydin@boun.edu.tr http://www.cmpe.boun.edu.tr/~ethem/i2ml3e



Rationale

□ Parameters θ not constant, but random variables with a prior, $p(\theta)$

□ Bayes' Rule:
$$p(\theta | X) = \frac{p(\theta)p(X | \theta)}{p(X)}$$

Generative Model



$$p(x', \mathcal{X}, \theta) = p(\theta) p(\mathcal{X}|\theta) p(x'|\theta)$$

$$p(x'|\mathcal{X}) = \frac{p(x',\mathcal{X})}{p(\mathcal{X})} = \frac{\int p(x',\mathcal{X},\theta)d\theta}{p(\mathcal{X})} = \frac{\int p(\theta)p(\mathcal{X}|\theta)p(x'|\theta)d\theta}{p(\mathcal{X})}$$
$$= \int p(x'|\theta)p(\theta|\mathcal{X})d\theta$$

Bayesian Approach

- 5
- 1. Prior $p(\theta)$ allows us to concentrate on region where θ is likely to lie, ignoring regions where it's unlikely
- 2. Instead of a single estimate with a single θ , we generate several estimates using several θ and average, weighted by how their probabilities
- Even if prior $p(\theta)$ is uninformative, (2) still helps.
- MAP estimator does not make use of (2):

 $\theta_{MAP} = \arg\max_{\theta} p(\theta | \mathcal{X})$

Bayesian Approach

$$p(x'|\mathcal{X}) = \int p(x'|\theta)p(\theta|\mathcal{X})d\theta$$

- In certain cases, it is easy to integrate
- Conjugate prior: Posterior has the same density as prior
- Sampling (Markov Chain Monte Carlo): Sample from the posterior and average
- Approximation: Approximate the posterior with a model easier to integrate
 - Laplace approximation: Use a Gaussian
 - Variational approximation: Split the multivariate density into a set of simpler densities using independencies

Estimating the Parameters of a Distribution: Discrete case

- $x_i^t = 1$ if in instance *t* is in state *i*, probability of state *i* is q_i
- Dirichlet prior, α_i are hyperparameters

Dirichlet(
$$\mathbf{q} \mid \boldsymbol{\alpha}$$
) = $\frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1)\cdots\Gamma(\alpha_K)} \prod_{i=1}^K q_i^{\alpha_i-1}$

$$p(X | \mathbf{q}) = \prod_{i=1}^{N} \prod_{j=1}^{K} q_{j}^{x_{j}^{t}}$$

• Posterior
$$p(\mathbf{q} \mid \boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0 + N)}{\Gamma(\alpha_1 + N_1) \cdots \Gamma(\alpha_K + N_K)} \prod_{i=1}^K q_i^{\alpha_i + N_i - 1}$$

$$=$$
 Dirichlet($\mathbf{q} \mid \boldsymbol{\alpha} + \mathbf{n}$)

Dirichlet is a conjugate prior

Sample likelihood

 \square With K=2, Dirichlet reduced to Beta

Estimating the Parameters of a Distribution: Continuous case

- p(x^t)~N(μ,σ²)
- Gaussian prior for μ , $p(\mu) \sim N(\mu_0, \sigma_0^2)$
- Posterior is also Gaussian $p(\mu|X) \sim N(\mu_{N'} \sigma_N^2)$ where





Gaussian: Prior on Variance

9

□ Let's define a prior (gamma) on precision $\lambda = 1/\sigma^2$ $p(\lambda) \sim \text{gamma}(a_0, b_0) = \frac{1}{\Gamma(a_0)} b_0^{a_0} \lambda^{a_0 - 1} \exp(-b_0 \lambda)$ $p(X|\lambda) = \prod_t \frac{\lambda^{1/2}}{\sqrt{2\pi}} \exp\left[-\frac{\lambda}{2} (x^t - \mu)^2\right]$ $= \lambda^{N/2} (2\pi)^{-N/2} \exp\left[-\frac{\lambda}{2} \sum_t (x^t - \mu)^2\right]$

 $p(\lambda|X) \propto p(X|\lambda)p(\lambda)$ ~ gamma(a_N, b_N)

$$a_N = a_0 + N/2 = \frac{\nu_0 + N}{2}$$
$$b_N = b_0 + \frac{N}{2}s^2 = \frac{\nu_0}{2}s_0^2 + \frac{N}{2}s^2$$

Joint Prior and Making a Prediction

 $p(\mu, \lambda) = p(\mu|\lambda)p(\lambda)$

 $p(\mu, \lambda | X) \sim \text{normal-gamma}(\mu_N, \kappa_N, a_N, b_N)$ where

 $\kappa_{N} = \kappa_{0} + N$ $\mu_{N} = \frac{\kappa_{0}\mu_{0} + Nm}{\kappa_{N}}$ $a_{N} = a_{0} + N/2$ $b_{N} = b_{0} + \frac{N}{2}s^{2} + \frac{\kappa_{0}N}{2\kappa_{N}}(m - \mu_{0})^{2}$ $p(x|X) = \iint p(x|\mu, \lambda)p(\mu, \lambda|X)d\mu d\lambda$ $\sim t_{2a_{N}}\left(\mu_{N}, \frac{b_{N}(\kappa_{N} + 1)}{a_{N}\kappa_{N}}\right)$

Multivariate Gaussian

11

 $p(\mathbf{x}) \sim \mathcal{N}_d(\boldsymbol{\mu}, \boldsymbol{\Lambda})$ $p(\boldsymbol{\mu}|\boldsymbol{\Lambda}) \sim \mathcal{N}_d(\boldsymbol{\mu}_0, (1/\kappa_0)\boldsymbol{\Lambda})$ $p(\boldsymbol{\Lambda}) \sim \text{Wishart}(\nu_0, \mathbf{V}_0)$ $p(\mu, \Lambda) = p(\mu|\Lambda)p(\Lambda)$ ~ normal-Wishart($\boldsymbol{\mu}_0, \kappa_0, \nu_0, \mathbf{V}_0$) $p(\mu, \Lambda | \mathcal{X}) \sim \text{normal-Wishart}(\mu_N, \kappa_N, \nu_N, \mathbf{V}_N)$ $\kappa_N = \kappa_0 + N$ $\boldsymbol{\mu}_N = \frac{\kappa_0 \boldsymbol{\mu}_0 + N \boldsymbol{m}}{\kappa_N}$ $v_N = v_0 + N$ $\mathbf{V}_N = \left(\mathbf{V}_0^{-1} + \mathbf{C} + \frac{\kappa_0 N}{\kappa_N} (\boldsymbol{m} - \boldsymbol{\mu}_0) (\boldsymbol{m} - \boldsymbol{\mu}_0)^T\right)^{-1}$ $p(\mathbf{x}|\mathcal{X}) = \int \int p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}) p(\boldsymbol{\mu}, \boldsymbol{\Lambda}|\mathcal{X}) d\boldsymbol{\mu} d\boldsymbol{\Lambda}$ ~ $t_{\nu_N-d+1}\left(\boldsymbol{\mu}_N, \frac{\kappa_N+1}{\kappa_N(\nu_N-d+1)}(\mathbf{V}_N)^{-1}\right)$

Estimating the Parameters of a Function: Regression

- $r = w^T x + \varepsilon$, $p(\varepsilon) \sim N(0, 1/\beta)$, and $p(r^t | x^t, w, \beta) \sim N(w^T x^t, 1/\beta)$
- Log likelihood $L(\mathbf{r} | \mathbf{X}, \mathbf{w}, \beta) = \log \prod_{t} p(r^{t} | \mathbf{x}^{t}, \mathbf{w}, \beta)$ $= -N \log(\sqrt{2\pi}) + N \log \beta - \frac{\beta}{2} \sum_{t} (r^{t} - \mathbf{w}^{T} \mathbf{x}^{t})$

ML solution $\mathbf{W}_{ML} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{r}$

- Gaussian conjugate prior: $p(w) \sim N(0, 1/\alpha)$
- Posterior: $p(\mathbf{w} \mid \mathbf{X}) \sim N(\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N)$ where

 $\boldsymbol{\mu}_{N} = \boldsymbol{\beta} \boldsymbol{\Sigma}_{N} \mathbf{X}^{\mathsf{T}} \mathbf{r}$ $\boldsymbol{\Sigma}_{N} = (\boldsymbol{\alpha} \mathbf{I} + \boldsymbol{\beta} \mathbf{X}^{\mathsf{T}} \mathbf{X})^{-1}$

Aka ridge regression/parameter shrinkage/ L2 regularization/weight decay



Prior on Noise Variance

14

 $p(\beta) \sim \text{gamma}(a_0, b_0) \qquad p(\boldsymbol{w}|\beta) \sim \mathcal{N}(\boldsymbol{\mu}_0, \beta \boldsymbol{\Sigma}_0)$

 $p(\mathbf{w}, \beta) = p(\beta)p(\mathbf{w}|\beta) \sim \text{normal-gamma}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0, \boldsymbol{a}_0, \boldsymbol{b}_0)$

 $p(\mathbf{w}, \beta | \mathbf{X}, \mathbf{r}) \sim \text{normal-gamma}(\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N, a_N, b_N)$

$$\Sigma_N = (\mathbf{X}^T \mathbf{X} + \Sigma_0)^{-1}$$

$$\mu_N = \Sigma_N (\mathbf{X}^T \mathbf{r} + \Sigma_0 \mu_0)$$

$$a_N = a_0 + N/2$$

$$b_N = b_0 + \frac{1}{2} (\mathbf{r}^T \mathbf{r} + \boldsymbol{\mu}_0^T \Sigma_0 \boldsymbol{\mu}_0 - \boldsymbol{\mu}_N^T \Sigma_N \boldsymbol{\mu}_N)$$

Markov Chain Monte Carlo (MCMC) sampling



Basis/Kernel Functions

16

• For new **x**', the estimate r' is calculated as

$$r' = (\mathbf{x}')^{T}$$

$$= \beta(\mathbf{x}')^{T} \Sigma_{N} \mathbf{X}^{T} \mathbf{r}$$

$$= \sum_{t} \beta(\mathbf{x}')^{T} \Sigma_{N} \mathbf{x}^{t} r^{t}$$
Dual representation

Linear kernel

• For any other $\phi(\mathbf{x})$, we can write $K(\mathbf{x}',\mathbf{x}) = \phi(\mathbf{x}')^T \phi(\mathbf{x})$

$$\mathbf{r'} = \sum_{t} \beta(\mathbf{x'})^{T} \Sigma_{N} \mathbf{x}^{t} \mathbf{r}^{t} \sum_{t} \beta K(\mathbf{x'}, \mathbf{x}^{t}) \mathbf{r}^{t}$$

Kernel Functions



What's in a Prior?

- Defining a prior is subjective
- Uninformative prior if no prior preference
- □ How high to go? Level I: $p(x|X) = \int p(x|\theta)p(\theta|X)d\theta$

Level II: $p(x|X) = \int p(x|\theta)p(\theta|X,\alpha)p(\alpha)d\theta d\alpha$

Empirical Bayes: Use one good α^*

Level II ML: $p(x|\mathcal{X}) = \int p(x|\theta)p(\theta|\mathcal{X}, \alpha^*)d\theta$

Bayesian Model Comparison

19

 Marginal likelihood of a model: p(X|M) = ∫ p(X|θ, M)p(θ|M)dθ

 Posterior probability of model given data: p(M|X) = p(X|M)p(M)/p(X) Bayes' factor: P(M_1|X)/P(M_0) = P(X|M_1)/P(M_0)/P(M_0)

Approximations:

BIC: $\log p(\mathcal{X}|\mathcal{M}) \approx \text{BIC} \equiv \log p(\mathcal{X}|\theta_{ML}, \mathcal{M}) - \frac{|\mathcal{M}|}{2} \log N$

AIC: AIC = log $p(\mathcal{X}|\theta_{ML}, \mathcal{M}) - |\mathcal{M}|$

Mixture Model

$$p(\mathbf{x}) = \sum_{i=1}^{k} P(G_i) p(\mathbf{x}|G_i)$$

$$p(\Phi) = p(\pi) \prod_i p(\mu_i, \Lambda_i)$$

$$= \text{Dirichlet}(\pi | \alpha) \prod_i \text{normal-Wishart}(\mu_0, \kappa_0, \nu_0, \mathbf{V}_0)$$

$$Q(\Phi | \Phi^l) = \sum_t \sum_i h_i^t \log \pi_i + \sum_t \sum_i h_i^t \log p_i(\mathbf{x}^t | \Phi^l) + \log p(\pi) + \sum_i \log p(\mu_i, \Lambda_i)$$

$$\pi_i^{l+1} = \frac{\alpha_i + N_i - 1}{\sum_i \alpha_i + N - k}$$

$$\mu_i^{l+1} = \left(\frac{\mathbf{V}_0^{-1} + \mathbf{C}_i + \mathbf{S}_i}{\nu_0 + N_i + d + 2}\right)^{-1}$$



Models in increasing complexity. A complex model can fit more datasets but is spread thin, a simple model can fit few datasets but has higher marginal likelihood where it does (MacKay 2003)



Nonparametric Bayes

- 22
- Model complexity can increase with more data (in practice up to N, potentially to infinity)
- Similar to k-NN and Parzen windows we saw before where training set is the parameters

Gaussian Processes

- Nonparametric model for supervised learning
- Assume Gaussian prior p(w)~N(0,1/α)
 y=Xw, where E[y]=0 and Cov(y)=K with K_{ij}= (xⁱ)^Txⁱ
 K is the covariance function, here linear
- With basis function $\phi(\mathbf{x})$, $\mathbf{K}_{ij} = (\phi(\mathbf{x}^i))^T \phi(\mathbf{x}^i)$ $r \sim N_N(\mathbf{0}, C_N)$ where $C_N = (1/\beta)\mathbf{I} + \mathbf{K}$
- With new **x**' added as \mathbf{x}_{N+1} , $r_{N+1} \sim N_{N+1}(0, C_{N+1})$

$$\mathbf{C}_{N+1} = \begin{bmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k} & \mathbf{c} \end{bmatrix}$$

where $\mathbf{k} = [K(\mathbf{x}',\mathbf{x}')_t]^T$ and $c=K(\mathbf{x}',\mathbf{x}')+1/\beta$. p(r'|x',X,r)~N($\mathbf{k}^T \mathbf{C}_{N-1}\mathbf{r}, c-\mathbf{k}^T \mathbf{C}_{N-1}\mathbf{k})$



Dirichlet Processes

- Nonparametric Bayesian approach for clustering
- Chinese restaurant process
- Customers arrive and either join one of the existing tables or start a new one, based on the table occupancies:

Join existing table *i* with $P(z_i = 1) = \frac{n_i}{\alpha + n - 1}, i = 1, ..., k$ Start new table with $P(z_{k+1} = 1) = \frac{\alpha}{\alpha + n - 1}$

Nonparametric Gaussian Mixture

26

Tables are Gaussian components and decisions based both on prior and also on input x:

Join component *i* with $P(z_i^t = 1) \propto$

Start new component with $P(z_{k+1}^t) \propto$

$$\frac{\frac{n_i}{\alpha + n - 1} p(\mathbf{x}^t | \mathcal{X}_i), i = 1, \dots, k}{\frac{\alpha}{\alpha + n - 1} p(\mathbf{x}^t)}$$

Latent Dirichlet Allocation

27

Bayesian feature extraction



Beta Processes

- Nonparametric Bayesian approach for feature extraction
- Matrix factorization:

$$\mathbf{X} = \mathbf{Z}\mathbf{A} \qquad \qquad z_j^t = \begin{cases} 1 & \text{with probability } \mu_j \\ 0 & \text{with probability } 1 - \mu_j \end{cases}$$

 $\mu_j \sim \text{beta}(\alpha, 1)$

- Nonparametric version: Allow *j* to increase with more data probabilistically
- Indian buffet process: Customer can take one of the existing dishes with prob µ_i or add a new dish to the buffet