

Beta Prior Distribution

$$P(\theta) = \frac{\theta^{\beta_H-1} (1-\theta)^{\beta_T-1}}{B(\beta_H, \beta_T)}$$

Conjugate Prior $P(\theta)$

specify which distribution it follows

covers only
material covered
through today
Conjugate Priors

$P(\theta|D)$ is in the same
distribution as $P(\theta)$.

ex $P(\theta) \sim \text{Beta}$, $P(\theta|D) \sim \text{Beta}$

Ex 1

is mean and mode same

$$\text{mean } \frac{\beta_H}{\beta_H + \beta_T} \quad \text{Mode } \frac{\beta_H - 1}{\beta_H + \beta_T - 2}$$

mean = mode iff $\beta_H = \beta_T$

$$p(D|\theta) \cdot P(\theta) = \underbrace{\theta^{\beta_H-1} (1-\theta)^{\beta_T-1} \cdot \frac{1}{B(\beta_H, \beta_T)}}_{\text{Prior}} \cdot \underbrace{\theta^{\alpha_H} (1-\theta)^{\alpha_T}}_{\text{Likelihood}}$$

$$\begin{aligned} \text{Posterior} &= P(\theta|D) = \frac{1}{C} p(D|\theta) P(\theta) \\ &= \frac{1}{C} \frac{1}{B(\beta_H, \beta_T)} \theta^{\alpha_H + \beta_H - 1} (1-\theta)^{\alpha_T + \beta_T - 1} \\ &\sim \text{Beta}(\alpha_H + \beta_H, \alpha_T + \beta_T) \\ &\text{// How does the } \frac{1}{C} \text{ factor in?} \end{aligned}$$

Mean = mode

$\Theta(\vec{D}) \sim \text{mean}$

$$\frac{\alpha_H + \beta_H}{\alpha_H + \beta_H + \alpha_T + \beta_T}$$

$$\text{MLE} \quad \frac{\alpha_H}{\alpha_H + \alpha_T}$$

Ex 2 MAP for Beta

$$P(\theta|D)$$

// How is MLE different from MAP

Prior: $\Theta \sim \text{Beta}(2, 3)$

① 2 heads 2 tails

$$\hat{\Theta}_{MLE} = \frac{2}{2+2} = \frac{1}{2}$$

$$\hat{\Theta}_{MAP} = \frac{4-1}{4+5-2} = \frac{3}{7}$$

$$P(\theta|D) \sim \text{Beta}(4, 5) = \text{Beta}(2+2, 3+2)$$

2 200 heads 200 tails

$$\hat{\Theta}_{MLE} = \frac{200}{200+200} = \frac{1}{2}$$

$$\hat{\Theta}_{MAP} = \frac{200 - 1 + 2}{202 + 203 - 2} = \frac{201}{403}$$

Ex

$$\text{posterior } \mathcal{B}(30, 40) \quad \text{Prior } \mathcal{B}(2, 3)$$

$$28 H \quad 37 T \quad B_H = 2 \quad B_T = 3$$

$$\text{Posterior: } \mathcal{B}(B_H + \alpha_H, B_T + \alpha_T)$$

$$30 = 2 + \alpha_H, 3 + \alpha_T = 40$$

$$\alpha_H = 28, \alpha_T = 37$$

Gaussians Properties Pg 19

How do we learn μ, σ

MLE for Gaussian

$$P(D|w, \sigma) = \frac{1}{\sigma \sqrt{2\pi}}^N \prod_{i=1}^N e^{-\frac{(x_i - w)^2}{2\sigma^2}}$$

$$\ln(P(D|w, \sigma)) = \ln \dots$$

$$= -N \ln \sigma \sqrt{2\pi} - \sum_{i=1}^N \frac{(x_i - w)^2}{2\sigma^2}$$

MLE for w (mean)

hold σ constant

$$\frac{\partial}{\partial w} \mathcal{L}(w, \sigma) = -N \ln \sigma \sqrt{2\pi} - \sum_{i=1}^N \frac{(x_i - w)^2}{2\sigma^2}$$

$$= -\sum_{i=1}^N \frac{(x_i - w)^2}{2\sigma^2}$$

$$= -\sum_{i=1}^N \frac{1}{2\sigma^2} (x_i - w)^2$$

$$= \cancel{\sum_{i=1}^N \frac{1}{2\sigma^2} - 2(x_i - w)} (x_i - w)$$

$$= \sum_{i=1}^N \frac{x_i - w}{\sigma^2} \cdot \sigma^2 = 0 \cdot \sigma^2$$

$$\sum_i (x_i - w) = 0$$

∴

$$\hat{w}_{MLE} = \frac{1}{N} \sum_i x_i$$

$E(\hat{w}_{MLE})$

$$\hat{w}_{MLE} = \frac{1}{N} \sum_i x_i = \underline{\sum_i (x_i - \bar{x}) = 0} \quad \text{if } \bar{x} \text{ is constant}$$

$$E(\hat{\mu}_{MLE}) = \frac{1}{n} \sum_i E(x_i) = \frac{1}{n} \sum_i \mu = \mu \quad \boxed{\frac{1}{n} \sum_i (x_i) = \mu}$$

= μ A true parameter of mean
 $\Rightarrow \hat{\mu}_{MLE}$ is unbiased

Properties of MLE Fisher info
 ↗ OPTIONAL ↗

MLE Variance p24

$$\begin{aligned} & \frac{\partial}{\partial \sigma} \left(-N \ln \sigma \sqrt{2\pi} - \frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2 \right) \\ &= + \frac{N}{\sigma} + \cancel{\frac{1}{2} \sum_i \frac{(x_i - \mu)^2}{\sigma^3}} - N \ln \sigma \sqrt{2\pi} - \cancel{2 \frac{1}{\sigma^2} \sum_i (x_i - \mu)^2} \\ &= - \frac{N}{\sigma \sqrt{2\pi}} + \cancel{2 \frac{1}{\sigma^2} \sum_i (x_i - \mu)^2} \\ & \hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_i (x_i - \mu)^2 = 0 = \frac{N}{\sigma^2} + \cancel{\frac{1}{\sigma^2} \sum_i (x_i - \mu)^2} \Rightarrow \frac{N - \sigma^2}{\sigma^2} = \cancel{\frac{1}{\sigma^2} \sum_i (x_i - \mu)^2} \\ & \Rightarrow \hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_i (x_i - \mu)^2 \end{aligned}$$

\Rightarrow USC MLE for σ since we do not have free sample variance

$$\begin{aligned} \hat{\sigma}_{MLE}^2 & \text{ is biased. Why? It uses an estimated } \hat{\mu} \\ \hat{\sigma}_{unbiased}^2 &= \frac{1}{N-1} \sum_i (x_i - \hat{\mu})^2 \\ &= \frac{1}{N} \sum_i (x_i - \mu)^2 \end{aligned}$$

check $E(\hat{\sigma}_{MLE}^2)$
 should be σ^2
 if not, then biased

Bayesian Learning of Gaussian Distribution

Conjugate Prior

Mean Gaussian Prior

Variance Wishart Distribution (\leftarrow not covered)

Prior for mean

$$P(\boldsymbol{\mu} | \gamma, \lambda) = \frac{1}{2\sqrt{2\pi}} e^{-\frac{(\boldsymbol{\mu} - \boldsymbol{\eta})^2}{2\lambda}}$$

γ, λ hyperparameters

MAP Gaussian Mean

$$P(\boldsymbol{\mu} | \gamma, \lambda) = \frac{1}{\lambda\sqrt{2\pi}} e^{-\frac{(\boldsymbol{\mu} - \boldsymbol{\eta})^2}{2\lambda^2}}$$

likelihood * prior

$$P(D | \boldsymbol{\mu}, \sigma^2) = \left(\frac{1}{\sigma\sqrt{2\pi}} \right)^N \prod_{i=1}^N e^{-\frac{(x_i - \boldsymbol{\mu})^2}{2\sigma^2}}$$

$$\frac{\partial}{\partial \boldsymbol{\mu}} \left[(\ln P(D | \boldsymbol{\mu}, \sigma^2)) + P(\boldsymbol{\mu} | \gamma, \lambda) \right] \xrightarrow{\text{argmax } \text{MLE} \Rightarrow \text{MAP for } \boldsymbol{\mu}}$$

$$\ln \frac{1}{\lambda\sqrt{2\pi}} - \frac{(\boldsymbol{\mu} - \boldsymbol{\eta})^2}{2\lambda^2} + N \ln \frac{1}{\sigma\sqrt{2\pi}} - \sum_{i=1}^N \frac{(x_i - \boldsymbol{\mu})^2}{2\sigma^2}$$

$$+ \frac{1}{2\lambda^2} \cancel{2(\boldsymbol{\mu} - \boldsymbol{\eta})} + \sum_{i=1}^N \frac{1}{2\sigma^2} \cancel{2(x_i - \boldsymbol{\mu})} = 0$$

$$- \frac{1}{\lambda^2} \boldsymbol{\mu} + \frac{\boldsymbol{\eta}}{\lambda^2} + \frac{1}{\sigma^2} \sum_i x_i - \frac{N}{\sigma^2} \boldsymbol{\mu} = 0$$

$$\hat{\boldsymbol{\mu}}_{\text{MAP}} = \begin{pmatrix} \frac{1}{\sigma^2} \sum_i x_i + \frac{\boldsymbol{\eta}}{\lambda^2} \\ \frac{N}{\sigma^2} + \frac{1}{\lambda^2} \end{pmatrix}$$

$$\bar{x} = \frac{1}{n} \sum x_i$$

Frequentist statistics
 Bayesian statistics
 probability as dynamics of belief

4:26 pm of the last day:
 we start on regression

Regression Minimize
 $w^* = E = \text{residual}$
 $= \arg \min_w (Hw - t)^T (Hw - t)$
 $\frac{\partial}{\partial w} = w^T H^T H w - 2t^T H w + t^T t$

$$2H^T H w - 2H^T t = 0$$

$$H^T H w = H^T t \quad \text{Pseudo inverse ?}$$

$$w^* = (H^T H)^{-1} H^T t$$

Maximize log-likelihood

$$\ln(P(D|w, \sigma)) = \ln\left(\frac{1}{\sigma^n} \prod_{i=1}^n e^{-\frac{(t_i - \sum w_i h_i(x_i))^2}{2\sigma^2}}\right)$$

assumed noise $\sim N(0, 1)$

$$\frac{\partial}{\partial w} \ln\left(\frac{1}{\sigma^n} \prod_{i=1}^n e^{-\frac{(t_i - \sum w_i h_i(x_i))^2}{2\sigma^2}}\right)$$

$$= \frac{2}{2\sigma} \sum_{i=1}^n -\frac{(t_i - \sum w_i h_i(x_i))^2}{2\sigma^2}$$

$$= \sum (t_i - \sum w_i h_i(x_i))^2 \leftarrow \begin{matrix} \text{residual} \\ \text{or} \\ \text{SSE} \end{matrix}$$

Least-squares Linear Regression for Gaussian

linear algebra slides for offline review

link for regression on BB next week

SVD not on First midterm

- skip Lecture 3 -
Start Lecture 4

Theorem Bayes classifier is optimal
 H_{Bayes}
 $\text{error}_{\text{true}}(h_{\text{Bayes}}) \leq \text{error}_{\text{true}}(h), \forall h(x)$ \leftarrow this is provable

$$\text{proof: } p(\text{error}) = \int_{\mathcal{X}} p(\text{error}(x)) p(x) dx$$

Bayes classifier

$$h(x) = \begin{cases} + & \text{if } p(Y=1) \geq 0.5 \\ - & \text{if } p(Y=1) < 0.5 \end{cases}$$

Bayes Rule

$$P(Y|X) = \frac{P(X|Y) P(Y)}{P(X)} \quad \begin{matrix} \text{likelihood} \\ \text{prior} \end{matrix} \quad \text{marginal}$$

How hard to learn optimal classifiers

$$x \in \mathbb{R}^n$$

Prior $P(Y)$ of K classes (Y_1, \dots, Y_K)

if $K=2$, 1 param needed

else need $K-1$ parameters

n features in m -space

$p(X|Y)$ where x composed of n features

$2^n \cdot K$ how does binary take effect?
 $m^n - 1$ possible features

DISCRETE

$$2^k - 1 = \text{due to binary choice}$$

$$(m^n - 1)^{K-1} \quad p(x=1 | y=1) = 1 - " = p(x=0)$$

for k classes
and n binary ($m=2$) features

prior $p(\gamma)$ $\binom{k-1}{n}$ parameters
for each γ , $p(x|y=k)$
 $\binom{m^n - 1}{m^n - 1}$ parameters
 $K(m^n - 1) + k-1$ parameters over all
exponential in terms of the number
of features.

Naive Bayes Assumption

Assume features are conditionally independent

$$P(E, A | P_a) = P(E | P_a) \cdot P(A | P_a)$$

factors class

new feature space of $p(x|y)$ // no effect on $p(\gamma)$
 n features (binary) no correlation

Kn factors $\theta \gamma$

$$\begin{aligned} y^* = h_{NB}(x) &= \operatorname{argmax}_y P(y) P(x_{1:n} | y) \\ &= \operatorname{argmax}_y P(y) \prod_i p(x_i | y) \end{aligned} \quad \begin{matrix} \text{via conditional} \\ \text{independence} \\ \text{assumption} \end{matrix}$$

if assumption holds, NB optimal
usually not true, but good enough

MLE for NB

Prior $P(Y=y)$

$$\sum_i y_i$$

Likelihood $P(X_i=x_i | Y_i=y_i)$

$$P(X_i=x_i | Y_i=y_i) = \frac{\# X_i=x_i \text{ and } Y_i=y_i}{\# Y_i=y_i}$$

No MAP on Exam