

Midterm Cheatsheet

Tuesday, October 19, 2021 8:49 PM

Bayes' Rule

- Using Bayes' posterior prob of class w_i given measurement x is given by:

$$P(w_i|x) = \frac{P(x|w_i)P(w_i)}{P(x)} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

Where $P(x) = \sum_{j=1}^2 P(x|w_j)P(w_j)$ & scale factor

Decide w_1 if $P(w_1|x) > P(w_2|x)$, otherwise w_2

Decide w_1 if $P(x|w_1)P(w_1) > P(x|w_2)P(w_2)$; ... w_2

Likelihood Ratio Test

$$\frac{P(x|w_1)}{P(x|w_2)} > \frac{\lambda_{12} - \lambda_{21}}{\lambda_{21} - \lambda_{12}} \frac{P(w_1)}{P(w_2)}$$

Let $\theta_2 = \frac{\lambda_{12} - \lambda_{21}}{\lambda_{21} - \lambda_{12}} \frac{P(w_1)}{P(w_2)}$, then

Decide w_1 if $\frac{P(x|w_1)}{P(x|w_2)} > \theta_2$
for zero-one loss: $\lambda = [0, 1]$, $\theta_2 = \frac{P(w_1)}{P(w_2)} = \theta_2$

Penalize more on misclassifying w_2 to w_1 ,
e.g. $\lambda = \begin{cases} 0 & w_1 \\ 10 & w_2 \end{cases}$, $\theta_2 = \frac{P(w_1)}{P(w_2)} = \theta_2$

$$\text{Gaussian: } f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$e^{\frac{(x-\mu)^2}{2\sigma^2}} > 1 \rightarrow \text{Do log to simplify}$$

$\Rightarrow x < ? \rightarrow \text{decide } w_1$
 $x > ? \rightarrow \text{decide } w_2$

Minimum Risk Decision: choose the action that minimize the conditional risk. ($R(a_i|x) = 1 - P(w_i|x)$)

$$\min_a R(a_i|x) = \max_a P(w_i|x)$$

Decide w_i if $P(w_i|x) > P(w_j|x)$ for $j \neq i$.

* Minimize prob error, minimize error rate, minimize risk, etc.

Logistic Regression

Consider problem of 2-class classification. Using discriminative approach, we tend to maximize a likelihood function defined thru the posterior.

$$\begin{aligned} p(C_1|\Phi) &= y(\Phi) = \sigma(\Phi) && \text{if } \sigma = \text{sigmoid} \\ p(C_2|\Phi) &= 1 - p(C_1|\Phi) && \text{if } \sigma = \frac{1}{1+e^{-x}} \end{aligned}$$

MLE for Logistic Regression

for dataset $\{(t_n, x_n)\}_{n=1}^N$, where $t_n \in \{0, 1\}$, $\Phi_n = \Phi(x_n)$, $n = 1, \dots, N$. Likelihood Function:

$$p(x=1) = \frac{1}{2} \left(\frac{1}{3} K\left(\frac{x-x_1}{3}\right) + \frac{1}{3} K\left(\frac{x-x_2}{3}\right) + \dots \right)$$

where: $t_n =$

$$\text{Get cross entropy: } E(w) = -\ln p(t|w) = -\sum_n t_n \ln y_n + (1-t_n) \ln(1-y_n)$$

Taking the gradient (using derivative of sigmoid):

$$\nabla E(w) = \sum_{n=1}^N (y_n - t_n) \Phi_n \quad \text{if } \Phi_n = \text{inputs}$$

Error function w/ penalty: $L_2(w) = \sum_{i=1}^d (y_i \ln f(x_i) + (1-y_i) \ln(1-f(x_i))) + \frac{\lambda}{2} \|w\|^2$

Sigmoid Derivative

$$\frac{d\sigma(x)}{dx} \approx 1 \approx 1$$

Perception

1) Loop over training data: x_1, x_2, \dots, x_N

2) In t step, if x_k is correctly classified, then $w^{(t+1)} = w^{(t)}$, otherwise:

$$w^{(t+1)} = \begin{cases} w^{(t)} + \alpha x_k & \text{if } x_k \text{ is in + class (+1)} \\ w^{(t)} - \alpha x_k & \text{if } x_k \text{ is in - class (-1)} \end{cases}$$

3) Convergence theorem: regardless of initial choice of weights, if 2 classes are linearly separable, there exists w such that:

$$\begin{cases} w^T x > 0 & \text{if } x \text{ in + class} \\ w^T x \leq 0 & \text{if } x \text{ in - class} \end{cases}$$

Note that: \leftarrow ground truth (+1 or -1)

$$\text{weights} \rightarrow \text{t inputs} \quad t_i(w^T x_i) > 0 \iff x_i \text{ correctly classified}$$

Parzen Window vs. KNN

Density Estimation: $p_n(x) = \frac{K_n}{V_n}$

Parzen Window: fix window size (e.g. $V_n = \frac{1}{n}$) &

$K_n = K_n(x)$ is a function of x :

$$p_n(x) = \frac{K_n(x)}{V_n} = \frac{K_n(x)}{\frac{1}{n}} = \frac{n}{V_n} K_n(x)$$

KNN: Fix # of samples inside window (e.g. $K_n = 5n$)
& $V_n = V_n(x)$ is a function of x : $\rightarrow K_n = 5$ points in V_n
 $p_n(x) = \frac{K_n}{V_n} = \frac{5}{V_n(x)} \frac{1}{n} \rightarrow V_n = \text{volume of } R_n$
 $n = \text{total } n$

$\rightarrow p_n(x) = n^{\text{th}} \text{ estimate for } p(x) = \text{density estimate}$

Parzen Window

Hypercube of volume $V_n = h^n$

Kernel Function: $K(u) = \begin{cases} 1, & |u_i| \leq 1/2, i=1, \dots, d \\ 0, & \text{otherwise} \end{cases}$

$K(u) = \text{unit hypercube centered at origin}$

$K\left(\frac{x-x_i}{h_n}\right)$ will be 1 if X_i lies in cube, 0 otherwise

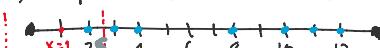
Points in hypercube: \rightarrow if $\in [0, 1]$

$$K_n = \sum_{i=1}^n K\left(\frac{x-x_i}{h_n}\right) \quad K_n = 1$$

Plug in $p_n(x)$ w/ V_n :

$$p_n(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{V_n} K\left(\frac{x-x_i}{h_n}\right)$$

Ex) 7 samples $D = \{2, 3, 4, 8, 10, 11, 12\}$, $h=3$, $x=1$



$h=3$ Only 1 point in hypercube

$$p(x=1) = \frac{1}{7} \left(\frac{1}{3} K\left(\frac{1-2}{3}\right) + \frac{1}{3} K\left(\frac{1-3}{3}\right) + \dots \right) = \frac{1}{7} \left(\frac{1}{3} \times 1 \right) = \frac{1}{21}$$

KNN Classification

$$p(C_i|x) = \frac{k_i}{n} \quad \text{choose class w/ highest } k_i$$

\rightarrow higher k = smoother, higher bias

MLE vs. Bayesian Estimation

MLE: trying to find value for params θ that maximize the likelihood. Not leaving for a distribution.

Bayesian: Fully calculate posterior to get distribution for the params. Bayesian is more complex but

$t_n = \text{labels}$ $\nabla E(w) = \sum_{i=1}^n (y_i - t_i) \Phi_i$ $\Phi_i = \text{inputs}$
 Error function w/ penalty: $J_2(w) = \sum_{i=1}^n \{-(y_i \log(\sigma(w^T x_i)) + (1-y_i) \log(1-\sigma(w^T x_i)))\} + \frac{\lambda}{2} \|w\|^2$
 Sigmoid Derivative

$$\begin{aligned}\frac{\partial \sigma(x)}{\partial x} &= \sigma(x)(1-\sigma(x)) \\ \nabla &= \sum_{i=1}^n -\frac{\partial}{\partial w} (y_i \log(\sigma(w^T x_i)) + (1-y_i) \log(1-\sigma(w^T x_i))) + \frac{\lambda}{2} \|w\|^2 \\ &= -\left[\frac{\partial}{\partial w} (\frac{y_i}{\sigma(w^T x_i)}) + \frac{1-y_i}{1-\sigma(w^T x_i)} \right] + \lambda w \\ &= -\left[\frac{y_i}{\sigma(w^T x_i)} - \frac{1-y_i}{1-\sigma(w^T x_i)} \right] \frac{\partial \sigma(w^T x_i)}{\partial w} + \lambda w \\ &= -\left[y_i(\frac{1-y_i}{1-\sigma(w^T x_i)}) - (1-y_i)(\frac{y_i}{\sigma(w^T x_i)}) \right] x_i + \lambda w \\ &= -(y_i(1-\sigma(w^T x_i)) - (1-y_i)\sigma(w^T x_i)) x_i + \lambda w \\ &= (y_i - \sigma(w^T x_i)) x_i + \lambda w \\ &= (\sigma(w^T x_i) - y_i) x_i + \lambda w \\ \Rightarrow &\boxed{\sum_{i=1}^n [\sigma(w^T x_i) - y_i] x_i + \lambda w}\end{aligned}$$

Discriminative: Directly assign x to a specific class.
Probabilistic Discriminative: Directly model posterior $p(t|x)$.
Probabilistic Generative: Model class-conditional + prior.

Likelihood + Prior
 Consider $t = w^T \Phi(x) + e$, $e \sim N(0, \beta^{-1})$. Assume noise precision β is known constant. Likelihood equals:

$$p(t|w) = \prod_{i=1}^n N(t_i | w^T \Phi(x_i), \beta^{-1})$$

Define Prior distribution on w:

$$p(w) = N(w | \mu_0, \Sigma_0)$$

Consider simplified Gaussian prior:
 Given $p(w|\alpha) = N(w | \alpha, \kappa^{-1} I)$ \Rightarrow only unknown

MLE: Trying to find value for params θ that maximize the likelihood. Not focusing for a distribution.
 Bayesian: Fully calculate posterior to get distribution for the params. Bayesian is more complex but more accurate since it has more information.

FLD vs. PCA

FLD: focusing on maximizing separability among known classes. It is supervised + reduces to line (ID).
 PCA: focusing on reducing dimensions to find features w/ highest variation. Unsupervised (no labels).

3 Types of Gradient Descent

Batch: Uses all data each iteration. Accurate + unbiased, but slow + easy to get trapped in local modes.

Mini-Batch: Use small batches each iteration. Faster + adds noise to prevent getting trapped. Noise may cause learning to go back & forth.

Stochastic: 1 piece of data each iteration. Most efficient but can be slow to converge.

Maximum Likelihood Estimator

$$\begin{aligned}f(x|\theta) &= \prod_{i=1}^n 2\theta x_i e^{-\theta x_i^2} \\ \ln(f(x|\theta)) &= \sum_{i=1}^n \ln(2\theta x_i e^{-\theta x_i^2}) \\ LL(\theta) &= \sum_{i=1}^n (\ln(2\theta x_i) - \ln(e^{-\theta x_i^2})) \\ LL(\theta) &= \sum_{i=1}^n (\ln(2\theta x_i) - \theta x_i^2) \\ \nabla_\theta LL(\theta) &= \sum_{i=1}^n \left(\frac{2x_i}{2\theta x_i} - x_i^2 \right) = \sum_{i=1}^n \left(\frac{1}{\theta} - x_i^2 \right) \\ \nabla_\theta LL(\theta) &= \frac{n}{\theta} - \sum_{i=1}^n x_i^2 \\ \text{Set } \nabla_\theta LL(\theta) &= 0 \rightarrow \frac{n}{\theta} - \sum_{i=1}^n x_i^2 = 0 \\ \frac{n}{\theta} &= \sum_{i=1}^n x_i^2 \rightarrow \theta = \frac{n}{\sum_{i=1}^n x_i^2} \quad \boxed{\text{MLE}}\end{aligned}$$

Likelihood for Uniform Distribution

$$\prod_{i=1}^n f(x_i; a, b) = \prod_{i=1}^n \frac{1}{(b-a)} = \frac{1}{(b-a)^n}$$

$$\log \prod_{i=1}^n f(x_i; a, b) = \log((b-a)^{-n}) = -n \log(b-a)$$

* Log Likelihood \uparrow

MLE: Largest possible a, smallest b

$$\text{Derivative wrt } a: \frac{n}{(b-a)}$$

$$\text{Derivative wrt } b: \frac{-n}{(b-a)}$$

$$a = \min \{x_i\}$$

$$b = \max \{x_i\}$$