

## Review sheet

Derivatives: Rules on my slides

Logs: Rules on my slides

Linear algebra:  $\mathbf{a}^T \mathbf{b}$  – also known as dot product of  $\mathbf{a}$  and  $\mathbf{b}$  or projection of  $\mathbf{b}$  onto  $\mathbf{a}$ , magnitude of  $\mathbf{a}$

Probability: Conditional, joint, marginal, Bayes rule  $P(A|B) = \frac{P(B|A)P(A)}{P(B)} \propto P(B|A)P(A)$

Gaussian distribution

Training and testing sets, training and testing error

Performing classification: take class that has highest probability (Bayes) or “probability” (logistic regression), or based on whether binary classifier output is positive or negative (SVM)

Generative classifiers:

Binary coin flips:

Likelihood:  $p^{|H|}(1-p)^{|T|}$

Prior: Beta distribution – role of  $\alpha$  and  $\beta$  parameters

MLE:  $p = \frac{|H|}{|H|+|T|}$

MAP:  $p = \frac{|H|+\alpha-1}{|H|+|T|+\alpha-1+\beta-1}$

Multi-variate binary classification:

$\theta_{ik}^j$  is probability  $i^{\text{th}}$  feature has value  $x_k$ , given the class is  $y_j$ :  $P(X_i=x_k | Y=y_j)$

$\theta^j$  is probability class is  $y_j$ :  $P(Y=y_j)$

MLE:  $P(X_i = x_k | Y = y_j; \theta_{ik}^j) = \frac{\#D(X_i=x_k \wedge Y=y_j)}{\#D(Y=y_j)}$

$$P(Y = y_j; \theta^j) = \frac{\#D(Y=y_j)}{|D|}$$

MAP:  $P(\theta_{ik}^j | X_i = x_k, Y = y_j) = \frac{\#D(X_i=x_k \wedge Y=y_j) + (\beta_j - 1)}{\#D(Y=y_j) + \sum_m (\beta_m - 1)}$

$$P(\theta^j | Y = y_j) = \frac{\#D(Y=y_j) + (\beta_j - 1)}{|D| + \sum_m (\beta_m - 1)}$$

Counting parameters:

Bayes with binary variables and two classes:  $2x(2^n-1)$

Naïve Bayes with binary variables and two classes:  $2xn$

Sigmoid:  $g(h) = \frac{1}{1+e^{-h}}$

Logistic regression: Classes 0 and 1 (or multi-class, see below)

$$p(y = 1|x; \theta) = g(\mathbf{w}^T \mathbf{x}) = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}}}$$

$$p(y = 0|x; \theta) = 1 - g(\mathbf{w}^T \mathbf{x}) = \frac{e^{-\mathbf{w}^T \mathbf{x}}}{1+e^{-\mathbf{w}^T \mathbf{x}}}$$

MLE update:

For each data point  $(y^i, \mathbf{x}^i)$ , update each weight  $w_j$ :

$$w_j \leftarrow w_j + \varepsilon x_j^i (y^i - g(\mathbf{w}^T \mathbf{x}^i))$$

MAP update:

$$L1: w_j \leftarrow w_j + \varepsilon \left[ x_j^i (y^i - g(\mathbf{w}^T \mathbf{x}^i)) - \frac{\text{sign}(w_j)}{\lambda N} \right]$$

$$L2: w_j \leftarrow w_j + \varepsilon \left[ x_j^i (y^i - g(\mathbf{w}^T \mathbf{x}^i)) - \frac{w_j}{\lambda N} \right]$$

Multi-class logistic regression:

For each data point  $(y^i, \mathbf{x}^i)$ , update each weight  $w_j$ :

$$w_{j,k} \leftarrow w_{j,k} + \varepsilon x_j^i (\delta(y^i - k) - p(y^i = k | \mathbf{x}^i; \theta))$$

$$\delta(a) = \begin{cases} 1 & \text{if } a = 0 \\ 0 & \text{otherwise} \end{cases}$$

Number of parameters: M classes, N features –  $(N+1) \times (M-1)$  regression parameters

Support vector machines: Classes -1 and +1

Maximize margin, minimize  $\mathbf{w}^T \mathbf{w}$  subject to  $\mathbf{w}^T \mathbf{x}^+ + b \geq 1$  and  $\mathbf{w}^T \mathbf{x}^- + b \leq -1$

Solution to max-margin problem:  $w = \sum_i \alpha^i \mathbf{x}^i y^i$      $\alpha^i \geq 0 \forall i$      $\sum_i \alpha^i y^i = 0$

Slack variables – penalizing errors

argmin<sub>w,b</sub>  $\mathbf{w}^T \mathbf{w} + C \sum_i \varepsilon_i$  subject to  $\mathbf{w}^T \mathbf{x}^+ + b \geq 1 - \varepsilon_i$  and  $\mathbf{w}^T \mathbf{x}^- + b \leq -1 + \varepsilon_i$   
 $\varepsilon_i \geq 0 \forall i$

Mapping functions: defining higher dimensions to achieve linear separation

$$\sum_i \alpha_i y_i \varphi(\mathbf{x}^i)^T \varphi(\mathbf{x}^j) + b$$

Represent dot product as kernel:  $K(\mathbf{x}^i, \mathbf{x}^j) = \varphi(\mathbf{x}^i)^T \varphi(\mathbf{x}^j)$

MLE vs. MAP

Number of parameters:

Deriving parameter estimates: log-probability, set derivative to 0 (our Bayesian examples so far)  
or change parameter in direction of derivative (our logistic regression examples so far)

Gaussian

MLE

MAP