

Gaussian NB

$$\hat{\mu}_{ik} = \frac{1}{\sum_j \delta(Y^j = y_k)} \sum_j X_i^j \delta(Y^j = y_k) \quad \hat{\sigma}_{ik}^2 = \frac{1}{\sum_j \delta(Y^j = y_k) - 1} \sum_j (X_i^j - \hat{\mu}_{ik})^2 \delta(Y^j = y_k) \quad \hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2$$

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

MAP $\hat{\theta} = \arg \max_{\theta} P(\theta | \mathcal{D})$ $\hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i$ MLE $\hat{\theta} = \arg \max_{\theta} P(\mathcal{D} | \theta)$

Derivatives

$$\begin{aligned} (cf)' &= c f'(x) & \frac{d}{dx}(a^x) &= a^x \ln(a) \\ (f \pm g)' &= f'(x) \pm g'(x) & \frac{d}{dx}(e^x) &= e^x \\ (fg)' &= f'g + fg' - \mathbf{Pr} & \frac{d}{dx}(\ln(x)) &= \frac{1}{x}, x > 0 \\ \left(\frac{f}{g}\right)' &= \frac{f'g - fg'}{g^2} - \mathbf{Qu} & \frac{d}{dx}(\ln|x|) &= \frac{1}{x}, x \neq 0 \\ & & \frac{d}{dx}(\log_a(x)) &= \frac{1}{x \ln a}, x > 0 \end{aligned}$$

SVM

$$\begin{aligned} g(\mathbf{x}) &= \mathbf{w}^T \mathbf{x} + b & g(\mathbf{x}) &= \sum_{i \in \text{SV}} \alpha_i K(\mathbf{x}_i, \mathbf{x}) + b \\ \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1) &= 0 & \mathbf{w} &= \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i = \sum_{i \in \text{SV}} \alpha_i y_i \mathbf{x}_i \\ \text{minimize } \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i & & \text{minimize } L_p(\mathbf{w}, b, \alpha_i) &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1) \\ \text{such that } & & \text{s.t. } & \alpha_i \geq 0 \\ y_i (\mathbf{w}^T \mathbf{x}_i + b) &\geq 1 - \xi_i & \text{maximize } & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \xi_i &\geq 0 & \text{s.t. } & \alpha_i \geq 0, \text{ and } \sum_{i=1}^n \alpha_i y_i = 0 \\ & & & 0 \leq \alpha_i \leq C \end{aligned}$$

Perceptron

AND $f(X_1, \dots, X_n) = X_1 \wedge \dots \wedge X_k \wedge \neg X_{k+1} \wedge \dots \wedge \neg X_n$ OR $g(X_1, \dots, X_n) = X_1 \vee \dots \vee X_k \vee \neg X_{k+1} \vee \dots \vee \neg X_n$

$w_0 = -k + 0.5$; $w_1 = \dots = w_k = 1$ $w_0 = n - k - 0.5$; $w_0 = n - 0.5$ $+1$ to a and -1 to b .
 $w_0 = -n + 0.5$; $w_{k+1} = \dots = w_n = -1$. $w_0 = \frac{a+b}{2}$ and $w_1 = \frac{a-b}{2}$.

Neural Net

- Initialize w_i 's and w_0 to random values.
- Until Convergence do
 - $\Delta w_0 = 0$
 - $\Delta w_i = 0$ for $i = 1$ to n
 - For each training example indexed by j do
 - Compute $o_j = w_0 + \sum_{i=1}^n w_i (x_i + x_i^{1.5})$
 - For each training example indexed by j do
 - $\Delta w_0 = \Delta w_0 + (y_j - o_j)^2$
 - $\Delta w_i = \Delta w_i + (y_j - o_j)^2 (x_{i,j} + x_{i,j}^{1.5})$ for $i = 1$ to n
 - $w_0 = w_0 + \eta \Delta w_0$
 - $w_i = w_i + \eta \Delta w_i$ for $i = 1$ to n

Backpropagation Algorithm

- Initialize all weights to small random numbers
 Until convergence, Do
 For each training example, Do
- Input it to network and compute network outputs
 - For each output unit k

$$\delta_k \leftarrow o_k(1 - o_k)(t_k - o_k)$$
 - For each hidden unit h

$$\delta_h \leftarrow o_h(1 - o_h) \sum_{k \in \text{outputs}} w_{h,k} \delta_k$$
 - Update each network weight $w_{i,j}$

$$w_{i,j} \leftarrow w_{i,j} + \Delta w_{i,j}$$

where $\Delta w_{i,j} = \eta \delta_j x_{i,j}$

Linear

$$\begin{aligned} \frac{\partial E}{\partial w_i} &= \sum_d (t_d - o_d)(-x_{i,d}) \\ \frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_d (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_d \frac{\partial}{\partial w_i} (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \\ &= \sum_d (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - \vec{w} \cdot \vec{x}_d) \end{aligned}$$

Sigmoid

$$\begin{aligned} \frac{\partial E}{\partial w_i} &= - \sum_{d \in D} (t_d - o_d) o_d (1 - o_d) x_{i,d} \\ \delta_h &\leftarrow o_h(1 - o_h) \sum_{k \in \text{outputs}} w_{h,k} \delta_k \end{aligned}$$

Minkowski distance

$$\begin{aligned} L_k(X_i, X_j) &= \left(\sum_{a=1}^d |x_{i,a} - x_{j,a}|^k \right)^{\frac{1}{k}} \\ L_{\infty}(X_i, X_j) &= \max_a |x_{i,a} - x_{j,a}| \end{aligned}$$

$\tanh \rightarrow f(z) = 1 - (f(z))^2$
 $\tanh(z)$ rescaled sigmoid $[-1, 1]$
 instead of $[0, 1]$

Linear reg

$$\begin{aligned} w_1 &= \frac{m \sum_{i=1}^m x_i y_i - \sum_{i=1}^m x_i \sum_{i=1}^m y_i}{m \sum_{i=1}^m x_i^2 - (\sum_{i=1}^m x_i)^2} & J(\mathbf{w}) &= \sum_{i=1}^m (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \\ w_0 &= \frac{\sum_{i=1}^m y_i - w_1 \sum_{i=1}^m x_i}{m} & \mathbf{w} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \end{aligned}$$

Logistic Reg

$$\begin{aligned} W &\leftarrow \arg \max_W \sum \ln P(Y^l | X^l, W) & P(Y = 1 | X) &= \frac{\exp(w_0 + \sum_{i=1}^n w_i X_i)}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)} \\ w_i &\leftarrow w_i + \eta \sum_l X_i^l (Y^l - \hat{P}(Y^l = 1 | X^l, W)) & P(Y = 0 | X) &= \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)} \\ l(W) &= \sum_l (Y^l (w_0 + \sum_i w_i X_i^l) - \ln(1 + \exp(w_0 + \sum_i w_i X_i^l))) & & \\ \frac{\partial l(W)}{\partial w_i} &= \sum_l X_i^l (Y^l - \hat{P}(Y^l = 1 | X^l, W)) & p(\mathbf{w}) &= \prod_i \frac{1}{\kappa \sqrt{2\pi}} e^{-\frac{w_i^2}{2\kappa^2}} \end{aligned}$$

$$w_i \leftarrow w_i + \eta \sum_l X_i^l (Y^l - \hat{P}(Y^l = 1 | X^l, W)) - \eta \lambda w_i$$

$$W \leftarrow \arg \max_W \sum_l \ln P(Y^l | X^l, W) - \frac{\lambda}{2} \|W\|^2$$

Pruning & regularization increases bias reduces variance
 Linear regressin: Loss +penalty (bcz minimizing)
 Log regression: Loss – penalty (bcz maximizing)
 Under zero mean, conditional independent variance assumption,
 LogR = Gaussian NB
 NB: Features independent given class ! assumption on $P(\mathbf{X}|\mathbf{Y})$
 LR: Functional form of $P(\mathbf{Y}|\mathbf{X})$, no assumption on $P(\mathbf{X}|\mathbf{Y})$
 therefore LR expected to outperform GNB
 Naïve Bayes needs $O(\log n)$ samples
 Logistic Regression needs $O(n)$ samples
 LR- conditional likelihood

HMM

Filtering: $P(\mathbf{X}_t|\mathbf{e}_{1:t})$

belief state—input to the decision process of a rational agent

Prediction: $P(\mathbf{X}_{t+k}|\mathbf{e}_{1:t})$ for $k > 0$

evaluation of possible action sequences;
 like filtering without the evidence

Smoothing: $P(\mathbf{X}_k|\mathbf{e}_{1:t})$ for $0 \leq k < t$

better estimate of past states, essential for learning

Most likely explanation: $\arg \max_{\mathbf{x}_{1:t}} P(\mathbf{x}_{1:t}|\mathbf{e}_{1:t})$

speech recognition, decoding with a noisy channel

$$\begin{aligned} P(\mathbf{X}_{t+1}|\mathbf{e}_{1:t+1}) &= P(\mathbf{X}_{t+1}|\mathbf{e}_{1:t}, \mathbf{e}_{t+1}) \\ &= \alpha P(\mathbf{e}_{t+1}|\mathbf{X}_{t+1}, \mathbf{e}_{1:t}) P(\mathbf{X}_{t+1}|\mathbf{e}_{1:t}) \\ &= \alpha P(\mathbf{e}_{t+1}|\mathbf{X}_{t+1}) P(\mathbf{X}_{t+1}|\mathbf{e}_{1:t}) \end{aligned}$$

Bayes Rule

Markov assumption

I.e., prediction + estimation. Prediction by summing out \mathbf{X}_t :

$$\begin{aligned} P(\mathbf{X}_{t+1}|\mathbf{e}_{1:t+1}) &= \alpha P(\mathbf{e}_{t+1}|\mathbf{X}_{t+1}) \sum_{\mathbf{x}_t} P(\mathbf{X}_{t+1}|\mathbf{x}_t, \mathbf{e}_{1:t}) P(\mathbf{x}_t|\mathbf{e}_{1:t}) \\ &= \alpha P(\mathbf{e}_{t+1}|\mathbf{X}_{t+1}) \sum_{\mathbf{x}_t} P(\mathbf{X}_{t+1}|\mathbf{x}_t) P(\mathbf{x}_t|\mathbf{e}_{1:t}) \end{aligned}$$

Forward-backward algorithm: cache forward messages along the way
 Time linear in t (polytree inference), space $O(t|f|)$

$\mathbf{f}_{1:t+1} = \text{FORWARD}(\mathbf{f}_{1:t}, \mathbf{e}_{t+1})$ where $\mathbf{f}_{1:t} = P(\mathbf{X}_t|\mathbf{e}_{1:t})$

Time and space **constant** (independent of t)

Divide evidence $\mathbf{e}_{1:t}$ into $\mathbf{e}_{1:k}$, $\mathbf{e}_{k+1:t}$:

$$\begin{aligned} P(\mathbf{X}_k|\mathbf{e}_{1:t}) &= P(\mathbf{X}_k|\mathbf{e}_{1:k}, \mathbf{e}_{k+1:t}) \\ &= \alpha P(\mathbf{X}_k|\mathbf{e}_{1:k}) P(\mathbf{e}_{k+1:t}|\mathbf{X}_k, \mathbf{e}_{1:k}) \\ &= \alpha P(\mathbf{X}_k|\mathbf{e}_{1:k}) P(\mathbf{e}_{k+1:t}|\mathbf{X}_k) \\ &= \alpha \mathbf{f}_{1:k} \mathbf{b}_{k+1:t} \end{aligned}$$

Bayes Rule

Markov assumption

Backward message computed by a backwards recursion:

$$\begin{aligned} P(\mathbf{e}_{k+1:t}|\mathbf{X}_k) &= \sum_{\mathbf{x}_{k+1}} P(\mathbf{e}_{k+1:t}|\mathbf{X}_k, \mathbf{x}_{k+1}) P(\mathbf{x}_{k+1}|\mathbf{X}_k) \\ &= \sum_{\mathbf{x}_{k+1}} P(\mathbf{e}_{k+1:t}|\mathbf{x}_{k+1}) P(\mathbf{x}_{k+1}|\mathbf{X}_k) \\ &= \sum_{\mathbf{x}_{k+1}} P(\mathbf{e}_{k+1}|\mathbf{x}_{k+1}) P(\mathbf{e}_{k+2:t}|\mathbf{x}_{k+1}) P(\mathbf{x}_{k+1}|\mathbf{X}_k) \end{aligned}$$

ADA BOOST

VC

• α_t :

- No errors: $\varepsilon_t=0 \rightarrow \alpha_t=\infty$
- All errors: $\varepsilon_t=1 \rightarrow \alpha_t=-\infty$
- Random: $\varepsilon_t=0.5 \rightarrow \alpha_t=0$

c belongs to C $m \geq \frac{1}{\epsilon} (\ln(1/\delta) + \ln(|H|))$

c not known $m \geq \frac{1}{2\epsilon^2} (\ln(1/\delta) + \ln(|H|))$

infinite |H| $m \geq \frac{1}{\epsilon} (4 \log_2(2/\delta) + 8VC(H) \log_2(13/\epsilon))$

Entropy

$$\text{GainRatio}(S, A) \equiv \frac{\text{Gain}(S, A)}{\text{SplitInformation}(S, A)}$$

$$H(V) = \sum_{v=0}^1 -P(H=v) \lg P(H=v).$$

$$\text{SplitInformation}(S, A) \equiv - \sum_{i=1}^c \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|}$$

$$\text{Gain}(S, A) = H(S) - \sum_{v \in \text{Values}(A)} P(v) H(S_v)$$

Beta Prior

Prior: $\text{Beta}(\beta_H, \beta_T)$

Data: α_H heads and α_T tails

$$\hat{\theta} = \arg \max_{\theta} P(\theta | \mathcal{D}) = \frac{\alpha_H + \beta_H - 1}{\alpha_H + \beta_H + \alpha_T + \beta_T - 2}$$

• Likelihood function: $P(\mathcal{D} | \theta) = \theta^{\alpha_H} (1 - \theta)^{\alpha_T}$

• Posterior: $P(\theta | \mathcal{D}) \propto P(\mathcal{D} | \theta) P(\theta)$

$$\begin{aligned} P(\theta | \mathcal{D}) &\propto \theta^{\alpha_H} (1 - \theta)^{\alpha_T} \theta^{\beta_H - 1} (1 - \theta)^{\beta_T - 1} \\ &= \theta^{\alpha_H + \beta_H - 1} (1 - \theta)^{\alpha_T + \beta_T - 1} \\ &= \text{Beta}(\alpha_H + \beta_H, \alpha_T + \beta_T) \end{aligned}$$

Mixture Model

• Given a dataset: $D = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$

• Mixture model: $\Theta = \{\alpha_1, \dots, \alpha_K, \theta_1, \dots, \theta_K\}$

$$p(\mathbf{x}|\Theta) = \sum_{k=1}^K \alpha_k p_k(\mathbf{x}|z_k, \theta_k)$$

The $p_k(\mathbf{x}|z_k, \theta_k)$ are mixture components, $1 \leq k \leq K$

$\mathbf{z} = (z_1, \dots, z_K)$ is a vector of K binary indicator variables

Note: only one of them equals 1 at any given point. Each point is assumed to be generated from exactly one mixture component!

$$\text{Mixture Weights. } \alpha_k = p(z_k) \quad \sum_{k=1}^K \alpha_k = 1.$$

the “membership weight” of data point \mathbf{x}_i in cluster k , given parameters Θ

$$w_{ik} = p(z_{ik} = 1 | \mathbf{x}_i, \Theta) = \frac{p_k(\mathbf{x}_i | z_k, \theta_k) \cdot \alpha_k}{\sum_{m=1}^K p_m(\mathbf{x}_i | z_m, \theta_m) \cdot \alpha_m}$$

Gaussian Mixture Models (GMMs)

$$p_k(\mathbf{x}|\theta_k) = \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{1/2}} e^{-\frac{1}{2}(\mathbf{x} - \mu_k)^t \Sigma_k^{-1} (\mathbf{x} - \mu_k)}$$

EM

Solution: As discussed in class, you can either update the sufficient statistics using the variable elimination algorithm for each example or complete each example. The complexity is the minimum of the two.

The complexity of completing the dataset is $O(2n \times 9)$ because there are two possible completions for each instance and there are nine functions which we have to multiply to compute the probability or weight for each example. Since the maximum CPT size is $O(\exp(3))$, we will need $O(9 \exp(3))$ to normalize the CPTs at the end. Thus the overall time complexity is $O(18n + 9 \exp(3))$.