

Evaluation:

- Problem - Compute Probability of observation sequence given a model
- Solution - **Forward Algorithm and Viterbi Algorithm**

Decoding:

- Problem - Find state sequence which maximizes probability of observation sequence
- Solution - **Viterbi Algorithm**

Training:

- Problem - Adjust model parameters to maximize probability of observed sequences
- Solution - **Forward-Backward Algorithm**

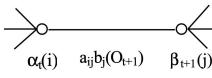
Baum-Welch Reestimation:

$$\bar{a}_{ij} = \frac{\text{expected number of trans from } S_i \text{ to } S_j}{\text{expected number of trans from } S_i} = \frac{\sum_{t=0}^{T-1} \xi_t(i, j)}{\sum_{t=0}^{T-1} N}$$

$$\bar{b}_j(k) = \frac{\text{expected number of times in state } j \text{ with symbol } k}{\text{expected number of times in state } j} = \frac{\sum_{t=0}^{T-1} \sum_{i=0}^N \xi_t(i, j)}{\sum_{t=0}^{T-1} N}$$

$\xi_t(i, j)$ = Probability of transitioning from S_i to S_j at time t given O

$$= P(q_t = S_i, q_{t+1} = S_j | O, \lambda) = \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{P(O | \lambda)}$$

**Convergence of F-B Algorithm:**

1. Initialize $\lambda = (\alpha, \beta)$
2. Compute α , β , and ξ
3. Estimate $\bar{\lambda} = (\bar{\alpha}, \bar{\beta})$ from ξ
4. Replace λ with $\bar{\lambda}$
5. If not converged go to 2

It can be shown that $P(O | \bar{\lambda}) > P(O | \lambda)$ unless $\bar{\lambda} = \lambda$

[Maximum Likelihood Estimation vs. Maximum a Posteriori]

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^N p(x^{(i)} | \theta) \quad \theta_{MAP} = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^N p(x^{(i)} | \theta) p(\theta)$$

MLE: As much prob mass as possible to things we have observed
MAP: Hallucinate data

MLE of exponential with $\lambda \exp(-\lambda x)$: MLE = $N / \Sigma x_i$

MLE of Bernoulli with $\phi_{MLE} = N_1 / (N_0 + N_1)$

MAP of Bernoulli-Beta(a,b): $\phi_{MAP} = (N_1+a-1) / (N_0+N_1+a+b-2)$

[Generative vs. Discriminative]

Generative: models joint $p(x, y)$, e.g. Naïve Bayes

Discriminative: models conditional $p(y | x)$, e.g. Logistic Reg

If model assumption is correct, NB is a more efficient learner (requires fewer samples). Otherwise, LR has lower asymptotic error.

Naïve Bayes:
Parameters are decoupled \rightarrow Closed form solution for MLE

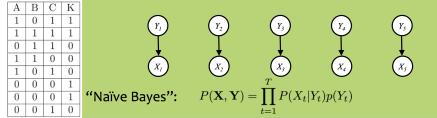
Bernoulli Naïve Bayes:
Parameters are probabilities \rightarrow Beta prior (usually) pushes probabilities away from zero / one extremes

Logistic Regression:
Parameters are not probabilities \rightarrow Gaussian prior encourages parameters to be close to zero (effectively pushes the probabilities away from zero / one extremes)

Logistic Regression:
Parameters are coupled \rightarrow No closed form solution - must use iterative optimization techniques instead

Naïve Bayes:
Features x are assumed to be conditionally independent given y . (i.e. Naïve Bayes Assumption)

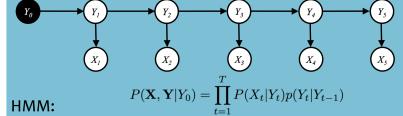
Logistic Regression:
No assumptions are made about the form of the features x . They can be dependent and correlated in any fashion.

[Naïve Bayes]

want $P(K=1 | A=1, B=1, C=0)?$

$$P(K=1 | A=1, B=1, C=0) = P(A=1, B=1, C=0 | K=1)P(K=1)$$

$$P(A=1, B=1, C=0 | K=1) = P(A=1 | K=1)P(B=1 | K=1)P(C=0 | K=1)$$

[Hidden Markov Model]

HMM:

For random variables X_1, X_2, X_3, X_4 :

$$\begin{aligned} P(X_1, X_2, X_3, X_4) &= P(X_1 | X_2, X_3, X_4) \\ &= P(X_2 | X_3, X_4) \\ &= P(X_3 | X_4) \\ &= P(X_4) \end{aligned}$$

3 problems for HMM:

1. Evaluation: Compute the probability of a given sequence of observations. $P(x) = \sum_y P(y, x)$
2. Decode: Find the most-likely sequence of hidden states, given observations. $y = \operatorname{argmax}_y P(y|x)$
3. Marginal: Compute the marginal distribution for a hidden state, given a sequence of observations. $P(y_t = k | x) = \sum_{y_1=y_t, y_{t+1}=k} P(y|x)$

Brutal force $O(K^T)$, Viterbi & F-B $O(TK^2)$

Derivation of Forward Algorithm

Definition: $\alpha_t(k) \triangleq p(x_1, \dots, x_t, y_t = k)$

$$\begin{aligned} \alpha_T(k_{END}) &= p(x_1, \dots, x_T, y_T = END) \\ &= p(x_1, \dots, x_T, y_T) p(y_T) \quad \text{[using "y_T" as shorthand for "y_T = END"]} \\ &= p(x_1, \dots, x_T, y_T) p(y_T | x_1, \dots, x_T) \\ &= p(x_1, \dots, x_T, y_T) \sum_{y_{T+1}} p(x_1, \dots, x_T, y_T, y_{T+1}) \\ &= p(x_1, \dots, x_T, y_T) \sum_{y_{T+1}} p(x_1, \dots, x_T, y_T) p(y_{T+1} | y_T) \\ &= p(x_1, \dots, x_T, y_T) \sum_{y_{T+1}} \alpha_{T+1}(y_{T+1}) p(y_{T+1} | y_T) \quad \text{[by def. of marginal]} \\ &= p(x_1, \dots, x_T, y_T) \sum_{y_{T+1}} \alpha_{T+1}(y_{T+1}) p(y_{T+1} | y_T) \quad \text{[by def. of joint]} \\ &= p(x_1, \dots, x_T, y_T) \sum_{y_{T+1}} \alpha_{T+1}(y_{T+1}) p(y_{T+1} | y_T) \quad \text{[by cond. indep. of HMM]} \\ &= p(x_1, \dots, x_T, y_T) \sum_{y_{T+1}} \alpha_{T+1}(y_{T+1}) p(y_{T+1} | y_T) \quad \text{[by def. of joint]} \\ &= p(x_1, \dots, x_T, y_T) \sum_{y_{T+1}} \alpha_{T+1}(y_{T+1}) p(y_{T+1} | y_T) \quad \text{[by cond. indep. of HMM]} \\ &= p(x_1, \dots, x_T, y_T) \sum_{y_{T+1}} \alpha_{T+1}(y_{T+1}) p(y_{T+1} | y_T) \quad \text{[by def. of joint]} \end{aligned}$$

Forward-Backward Algorithm

Define: $\alpha_t(k) \triangleq p(x_1, \dots, x_t, y_t = k)$ Assume $y_0 = \text{START}$
 $\beta_t(k) \triangleq p(x_{t+1}, \dots, x_T | y_t = k)$ $y_T = \text{END}$

- ① Initialize $\alpha_0(\text{START}) = 1$ $\alpha_0(k) = 0 \forall k \neq \text{START}$
 $\beta_{T+1}(\text{END}) = 1$ $\beta_{T+1}(k) = 0 \forall k \neq \text{END}$
- ② For $t = 1, \dots, T$:
 - For $k = 1, \dots, K$:

$$\alpha_t(k) = p(x_t | y_t = k) \sum_{j=1}^K \alpha_{t-1}(j) p(y_t = k | y_{t-1} = j)$$
 - For $t = T, \dots, 1$:
 - For $k = 1, \dots, K$:

$$\beta_t(k) = \sum_{j=1}^K p(x_{t+1} | y_{t+1} = j) \beta_{t+1}(j) p(y_{t+1} = j | y_t = k)$$
- ④ Compute $p(\vec{x}) = \alpha_{T+1}(\text{END})$ [Evaluation]
- ⑤ Compute $p(y_t = k | \vec{x}) = \alpha_t(k) \beta_t(k)$ [Marginals]

Viterbi Algorithm

Define: $\omega_t(k) \triangleq \max_{y_1, \dots, y_{t-1}} p(x_1, \dots, x_t, y_t = k)$

"backpointer" $\rightarrow b_t(k) \triangleq \operatorname{argmax}_{y_1, \dots, y_{t-1}} p(x_1, \dots, x_t, y_1, \dots, y_{t-1}, y_t = k)$

- Assume $y_0 = \text{START}$
- ① Initialize $\omega_0(\text{START}) = 1$ $\omega_0(k) = 0 \forall k \neq \text{START}$
 - ② For $t = 1, \dots, T$:
 - For $k = 1, \dots, K$:

$$\omega_t(k) = \max_{j \in \{1, \dots, K\}} p(x_t | y_t = k) \omega_{t-1}(j) p(y_t = k | y_{t-1} = j)$$
 - For $t = T, \dots, 1$:
 - For $k = 1, \dots, K$:

$$b_t(k) = \operatorname{argmax}_{j \in \{1, \dots, K\}} p(x_t | y_t = k) \omega_{t-1}(j) p(y_t = k | y_{t-1} = j)$$
 - ③ Compute Most Probable Assignment

$$\hat{y} = \operatorname{argmax}_{b_{T+1}(\text{END})}$$

$$\hat{y} = \operatorname{argmax}_{t=1, \dots, T} b_t(\hat{y}_{t-1})$$

follow the backpointer

Short Comings:

- HMM models capture dependences between each state and only its corresponding observation
- Mismatch between learning objective function $P(X, Y)$ and prediction objective function $P(Y|X)$

Minimum Bayes Risk Decoding: (0-1 loss == viterbi)

$$\begin{aligned} h_{\theta}(x) &= \operatorname{argmin}_{\hat{y}} \mathbb{E}_{y \sim p_{\theta}(y | x)} [\ell(\hat{y}, y)] \\ &= \operatorname{argmin}_{\hat{y}} \sum_{y} p_{\theta}(y | x) \ell(\hat{y}, y) \end{aligned}$$

Direct Graphical Models (Bayes Net)

$$\begin{aligned} p(X_1, X_2, X_3, X_4, X_5) &= p(X_5 | X_1, X_2, X_3, X_4) \\ &= p(X_5 | X_3) p(X_4 | X_2, X_3) \\ &= p(X_3) p(X_2 | X_1) p(X_1) \end{aligned}$$

In order for a Bayesian network to model a probability distribution, the following must be true:

Each variable is conditionally independent of all its non-descendants in the graph given the value of all its parents.

D-Separation:

If variables X and Z are **d-separated** given a set of variables E
Then X and Z are **conditionally independent** given the set E

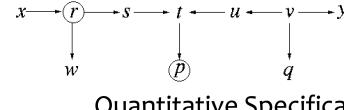
Definition #1:

Variables X and Z are **d-separated** given a set of evidence variables E iff every path from X to Z is "blocked".

A path is "blocked" whenever:

1. $\exists Y$ on path s.t. $Y \in E$ and Y is a "common parent"
2. $\exists Y$ on path s.t. $Y \in E$ and Y is in a "cascade"
3. $\exists Y$ on path s.t. $\{Y, \text{descendants}(Y)\} \notin E$ and Y is in a "v-structure"

Rule 3: If a collider is a member of the conditioning set Z , then it no longer blocks any path that traces this collider.

**Quantitative Specification**

Example: Conditional probability tables (CPTs) for discrete random variables

a^0	0.75
a^1	0.25

$$P(a, b, c, d) = P(a)P(b|a)P(c|a, b)P(d|c)$$

A	B
C	D
a^0	$a^0 b^0$
c^0	$a^0 b^0$
c^1	$a^1 b^0$
d^0	$a^1 b^1$
d^1	$a^1 b^1$

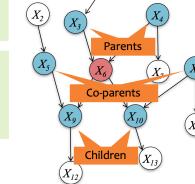
Markov Blanket

Def: the co-parents of a node are the parents of its children

Def: the Markov Blanket of a node is the set containing the node's parents, children, and co-parents.

Theorem: a node is conditionally independent of every other node in the graph given its Markov blanket

Example: The Markov Blanket of X_4 is $\{X_3, X_4, X_5, X_6, X_7, X_8, X_9, X_{10}\}$



Learning a Bayes Net: Learning this fully observed Bayesian Network is equivalent to learning five (small / simple) independent networks from the same data

[Reinforcement Learning]**Fixed point iteration:**

- Given objective function: $J(x) = \frac{x^3}{3} + \frac{3}{2}x^2 + 2x$
- Compute derivative, set to zero (call this function f): $dJ/dx = f(x) = x^2 - 3x + 2 = 0$
- Rearrange the equation s.t. one of parameters appears on the LHS.
- Initialize the parameters.
- For $i \in \{1, \dots, K\}$, update each parameter and increment i .
- Repeat #5 until convergence

Value iteration and **Policy iteration** are two more classic approaches to this problem. But essentially both are **dynamic programming**.

Q-learning is a more recent approaches to this problem. Essentially it is a **temporal-difference method**.

Algorithm 1 Value Iteration

- procedure **VALUEITERATION**($R(s, a)$ reward function, $p(\cdot | s, a)$ transition probabilities)
 - Initialize value function $V(s) = 0$ or randomly
 - while not converged do
 - for $s \in \mathcal{S}$ do
 - for $a \in \mathcal{A}$ do

$$Q(s, a) = R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a)V(s')$$
 - $V(s) = \max_a Q(s, a)$
- Let $\pi(s) = \operatorname{argmax}_a Q(s, a)$, $\forall s$
- return π

$O(|\mathcal{A}| |\mathcal{S}|^2)$ calculations per iteration, $O(|\mathcal{A}| |\mathcal{S}|)$ steps to converge (number of all possible policies)

Algorithm 2 Policy Iteration

- procedure **POLICYITERATION**($R(s, a)$ reward function, $p(\cdot | s, a)$ transition probabilities)
 - Initialize policy π randomly
 - while not converged do
 - Solve Bellman equations for fixed policy π

$$V^{\pi}(s) = R(s, \pi(s)) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, \pi(s))V^{\pi}(s')$$
 - Improve policy π using new value function

$$\pi(s) = \operatorname{argmax}_a R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a)V^{\pi}(s')$$
- return π

$O(|A| |S|^2 + |S|^3)$ calculations per iteration, converge faster.

Algorithm Q-Learning

$$V^*(s) = \max_a Q^*(s, a), \quad \pi^*(s) = \arg \max_a Q^*(s, a)$$

$$Q^*(s, a) = R(s, a) + \gamma \sum_{s'} p(s'|s, a) \max_{a'} Q^*(s', a')$$

$$1. \text{ Initialize } Q(s, a) = 0 \quad \forall s, a$$

2. Do forever

1. Select action a and execute

2. receive reward $r = R(s, a)$

[Deterministic]

3. observe new state $s' = \delta(s, a)$

4. update table entry $Q(s, a) = r + \gamma \max_{a'} Q(s', a')$

[Nondeterministic]

3. observe new state $s' \sim p(\cdot | s, a)$

4. update table entry $Q(s, a) = (1 - \alpha_n)Q(s, a) + \alpha_n(r + \gamma \max_{a'} Q(s', a'))$, where $\alpha_n = 1/(1 + K_n(s, a))$

$$\alpha_n = \frac{1}{1 + \# \text{visits at time stamp } n \text{ at } (s, a) \text{ pair}}$$

Epsilon Greedy:

w.p. $1 - \epsilon$, select $a = \max_{a'} Q(s, a')$,

w.p. ϵ , select $a = \text{random}$.

Support Vector Machine

Maximize “margin”

(There exists a \mathbf{w} and b , s.t. $y^{(i)}(\mathbf{w}^T \mathbf{x} + b) \geq 1$)

margin/width = $d_+ + d_-$

$$d_+ = \|\mathbf{w}\|_2 + b, \quad d_- = -\|\mathbf{w}\|_2 - b$$

Constraint is tight for x_+ and x_- , which means:

$$\mathbf{w}^T \mathbf{x}_+ = 1, \quad \mathbf{w}^T \mathbf{x}_- = -1,$$

minimize margin/width = $2 / \|\mathbf{w}\|_2$

maximize $\|\mathbf{w}\|_2 / 2$

Outliers doesn't matter for SVM but does for LR

Hard-margin SVM (Primal)

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1, \quad \forall i = 1, \dots, N$$

Hard-margin SVM (Lagrangian Dual)

$$\max_{\alpha} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y^{(i)} y^{(j)} \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}$$

$$\text{s.t. } \alpha_i \geq 0, \quad \forall i = 1, \dots, N$$

$$\sum_{i=1}^N \alpha_i y^{(i)} = 0$$



Hard-margin SVM (Primal)

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1, \quad \forall i = 1, \dots, N$$

$$\max_{\alpha} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y^{(i)} y^{(j)} \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}$$

$$\text{s.t. } \alpha_i \geq 0, \quad \forall i = 1, \dots, N$$

$$\sum_{i=1}^N \alpha_i y^{(i)} = 0$$

Soft-margin SVM (Primal)

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \left(\sum_{i=1}^N e_i \right)$$

$$\text{s.t. } y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - e_i, \quad \forall i = 1, \dots, N$$

$$\max_{\alpha} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y^{(i)} y^{(j)} \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}$$

$$\text{s.t. } 0 \leq \alpha_i \leq C, \quad \forall i = 1, \dots, N$$

$$e_i \geq 0, \quad \forall i = 1, \dots, N$$

$$\sum_{i=1}^N \alpha_i y^{(i)} = 0$$

one-vs-rest:

- build K binary classifiers
- train the k^{th} classifier to predict whether an instance has label k or something else
- predict the class with largest score

one-vs-one:

- build $\binom{K}{2}$ binary classifiers
- train one classifier for distinguishing between each pair of labels
- predict the class with the most “votes” from any given classifier

$$L(w, \alpha, \beta) = \frac{1}{2} \omega^T \omega - \frac{N}{2} \sum_{i=1}^N \alpha_i \zeta_i [y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)}) - 1]$$

$$\begin{aligned} \text{Primal: } & \min_{\mathbf{w}, b} \max_{\alpha, \beta} L(w, \alpha, \beta) \\ \text{Dual: } & \max_{\alpha, \beta} \min_{\mathbf{w}, b} L(w, \alpha, \beta) \end{aligned}$$

Lagrangian Multipliers

$$\begin{aligned} \min_{\mathbf{w}} & f(\mathbf{w}) \\ \text{s.t.} & g_i(\mathbf{w}) \leq 0, i = 1, \dots, k \\ & h_i(\mathbf{w}) = 0, i = 1, \dots, l \end{aligned}$$

This is also called the primal problem. We introduce new variables $\alpha_i \geq 0$ and β_i (no need to be non-negative) called Lagrange multipliers and study the generalized Lagrangian defined by

$$L(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w)$$

$$\text{Solve } \frac{\partial L(w, \alpha, \beta)}{\partial w} = 0, \quad \frac{\partial L(w, \alpha, \beta)}{\partial \alpha} = 0, \quad \frac{\partial L(w, \alpha, \beta)}{\partial \beta} = 0.$$

s.t. $\lambda \geq 0, g(\mathbf{w}) \leq 0$ (or C).

Matrix Factorization

SGD for UMF:

While not converged:

① Sample (i, j) from \mathbb{Z} uniformly at random

② Compute $e_{ij} = r_{ij} - \bar{v}_i^T \bar{v}_j$

③ Update $\bar{v}_i \leftarrow \bar{v}_i - \gamma \bar{v}_i^T \bar{J}_j(v_j)$

$\bar{v}_j \leftarrow \bar{v}_j - \gamma \bar{v}_j^T \bar{J}_i(v_i)$ W/ Regularization

$$\bar{J}_j(v_j) = \frac{1}{2} (r_{ij} - \bar{v}_i^T \bar{v}_j)^2 + \lambda (\|v_i\|_2^2 + \|v_j\|_2^2)$$

$$\bar{v}_i^T \bar{J}_j(v_j) = -e_{ij} \bar{v}_i^T + \lambda \bar{v}_i$$

$$\text{where } e_{ij} = r_{ij} - \bar{v}_i^T \bar{v}_j$$

Principle Component Analysis

Assume data is 0-centered and variance along each axis is 1.

Equivalence of Maximizing Variance and Minimizing Reconstruction Error

Claim: Minimizing the reconstruction error is equivalent to maximizing the variance.

Proof: First, note that:

$$\|\mathbf{x}^{(i)} - (\mathbf{v}^T \mathbf{x}^{(i)}) \mathbf{v}\|^2 = \|\mathbf{x}^{(i)}\|^2 - (\mathbf{v}^T \mathbf{x}^{(i)})^2 \quad (1)$$

since $\mathbf{v}^T \mathbf{v} = \|\mathbf{v}\|^2 = 1$.

Substituting into the minimization problem, and removing the extraneous terms, we obtain the maximization problem.

$$\mathbf{v}^* = \operatorname{argmax}_{\mathbf{v}: \|\mathbf{v}\|^2=1} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}^{(i)} - (\mathbf{v}^T \mathbf{x}^{(i)}) \mathbf{v}\|^2 \quad (2)$$

$$= \operatorname{argmin}_{\mathbf{v}: \|\mathbf{v}\|^2=1} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}^{(i)}\|^2 - (\mathbf{v}^T \mathbf{x}^{(i)})^2 \quad (3)$$

$$= \operatorname{argmax}_{\mathbf{v}: \|\mathbf{v}\|^2=1} \frac{1}{N} \sum_{i=1}^N (\mathbf{v}^T \mathbf{x}^{(i)})^2 \quad (4)$$

$$= \operatorname{argmax}_{\mathbf{v}: \|\mathbf{v}\|^2=1} \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)} \cdot \mathbf{x}^{(i)} \quad (5)$$

V is first K cols of V in SVD of X = USV^T

Ensemble Methods

Weighted Majority:

Suppose we have a pool of T binary classifiers $\mathcal{A} = \{h_1, \dots, h_T\}$ where $h_t : \mathbb{R}^M \rightarrow \{+1, -1\}$. Let α_t be the weight for classifier h_t .

Algorithm 1 Weighted Majority Algorithm

```

1: procedure WEIGHTEDMAJORITY( $\mathcal{A}, \beta$ )
2:   Initialize classifier weights  $\alpha_t = 1, \forall t \in \{1, \dots, T\}$ 
3:   for each training example  $(x, y)$  do
4:     Predict majority vote class (splitting ties randomly)

```

$$\hat{h}(x) = \operatorname{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$$

```

5:     if a mistake is made  $\hat{h}(x) \neq y$  then
6:       for each classifier  $t \in \{1, \dots, T\}$  do
7:         if  $h_t(x) \neq y$ , then  $\alpha_t \leftarrow \beta \alpha_t$ 

```

Theorem 0.1 (Littlestone & Warmuth, 1994). If the Weighted Majority Algorithm is applied to a pool \mathcal{A} of classifiers, and if each algorithm makes at most m mistakes on the sequence of examples, then the total number of mistakes is upper bounded by $2.4(\log |\mathcal{A}| + m)$.

AdaBoost

Given: $(x_1, y_1), \dots, (x_m, y_m)$ where $x_i \in X, y_i \in Y = \{-1, +1\}$

Initialize $D_1(i) = 1/m$.

For $t = 1, \dots, T$:

• Train weak learner using distribution D_t .

• Get weak hypothesis $h_t : X \rightarrow \{-1, +1\}$ with error

$$\epsilon_t = \Pr_{x \sim D_t} [h_t(x) \neq y_t].$$

• Choose $\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$.

• Update:

$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} \times \begin{cases} e^{-\alpha_t} & \text{if } h_t(x_i) = y_i \\ e^{\alpha_t} & \text{if } h_t(x_i) \neq y_i \\ \frac{D_t(i) \exp(-\alpha_t) h_t(x_i)}{Z_t} & \text{otherwise} \end{cases}$$

where Z_t is a normalization factor (chosen so that D_{t+1} will be a distribution).

Output the final hypothesis:

$$H(x) = \operatorname{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right).$$

Weak learners. AdaBoost never overfits.

[K]-Nearest Neighbor

2.1 Euclidean distance: $g(\mathbf{u}, \mathbf{v}) = \sqrt{\sum_{m=1}^M (u_m - v_m)^2}$

Manhattan distance: $g(\mathbf{u}, \mathbf{v}) = \sqrt{\sum_{m=1}^M |u_m - v_m|}$

2.2 Questions on KNN

1. How to handle even values of k ?

1. keep adding points 2. distance-weighted 3. remove furthest point...

2. What is the inductive bias of KNN?

1. Similar points should have similar labels

2. All dimensions are created equally

3. Computational Efficiency of KNN?

o Suppose we have N training examples, and each one has M features.

o When training, Naive ($k = 1$) is $O(1)$, k-d Tree is $O(MN \log N)$

o When testing, Naive ($k = 1$) is $O(MN)$, k-d Tree is $O(2^M \log N)$

4. Theoretical Guarantees of KNN?

o $\text{error}_{true}(h) < 2 \times \text{Bayes error rate}$ (Cover & Hart, 1967)

2.3 Cross Validation

Leave out different folds as validation sets, and get different validation error curves. Then we can average the curves and get a smoother validation error curve, and finally choose model with lowest validation error.

Perceptron

3.1 Definition

$h(\mathbf{x}) = \operatorname{sign}(\mathbf{w}^T \mathbf{x} + b)$ where $\operatorname{sign}(x) = 1$ if $x \geq 0$ else -1

Note that $0 = \mathbf{w}^T \mathbf{x} + b$ is orthogonal to \mathbf{w}

3.2 (Online) Perceptron Algorithm

Initialize parameters: weights $\mathbf{w} = [w_1, w_2, \dots, w_M]^T = \mathbf{0}$, bias $b = 0$.

for $i = 1, 2, 3, \dots$

1. Receive next example $(\mathbf{x}^{(i)}, y^{(i)})$

2. Predict $\hat{y} = \operatorname{sign}(\mathbf{w}^T \mathbf{x}^{(i)} + b)$

3. Switch on cases

o If false negative, i.e. $\hat{y} \neq y^{(i)} = 1$, $\mathbf{w} = \mathbf{w} + \mathbf{x}^{(i)}, b = b + 1$

o If false positive, i.e. $\hat{y} \neq y^{(i)} = -1$, $\mathbf{w} = \mathbf{w} - \mathbf{x}^{(i)}, b = b - 1$

Note: the hypothesis space of perceptron algorithm is all possible linear decision boundaries: $\mathcal{H} = \{h(\cdot) \mid \exists \mathbf{w} \in \mathbb{R}^M, b \in R, \text{ s.t. } h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b\}$

Why this algorithm works?

• for $|\mathbf{w}|_2 = 1$, $\mathbf{w}^T \mathbf{x}$ is the length of vector project of \mathbf{x} onto \mathbf{w}

• wrongly predicted \mathbf{w} corrects \mathbf{w}

• \mathbf{b} shifts the zero points where \mathbf{w} starts

Note: the inductive bias of perceptron algorithm is that

• decision boundary should be linear

• prefers to correct most recent mistakes

3.3 Batch Perceptron Algorithm

```

1: procedure PERCEPTRON( $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$ )
2:    $\mathbf{w} \leftarrow \mathbf{0}$  Initialize parameters
3:   while not converged do
4:     for  $i \in \{1, 2, \dots, N\}$  do For each example
5:        $\hat{y} \leftarrow \operatorname{sign}(\mathbf{w}^T \mathbf{x}^{(i)})$  Predict
6:       if  $\hat{y} \neq y^{(i)}$  then If mistake
7:          $\mathbf{w} \leftarrow \mathbf{w} + y^{(i)} \mathbf{x}^{(i)}$  Update parameters
8:     return  $\mathbf{w}$ 

```

Note: It cannot converge on not linearly separable data

3.4 Perceptron Mistake Bound:

With margin γ and all points inside a ball of radius R , it makes $\leq (R / \gamma)^2$ mistakes

Linear Regression and Optimization

4.1 Linear Regression as Function Approximation

(1) Assume \mathbf{D} is generated as $\mathbf{x}^{(i)} \sim p^*(\mathbf{x}), y^{(i)} \sim c^*(\mathbf{x}^{(i)})$, where p^* and c^* are unknown

(2) Choose $\mathcal{H} = \{h(\cdot) | h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}, \mathbf{w} \in \mathbb{R}^{M+1}, x_0 = 1\}$

(3) Choose a objective function with a goal of minimizing mean squared

$$\mathcal{J}(\mathbf{D}, \mathbf{w}) = \frac{1}{N} \sum_{i=1}^N e_i^2 = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - (\mathbf{w}^T \mathbf{x}^{(i)} + b))^2$$

(4) Solve optimization problem $\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \mathcal{J}(\mathbf{D}, \mathbf{w})$

(5) Predict, given next \mathbf{x} , $y = h_{\hat{\mathbf{w}}}(\mathbf{x}) = \hat{\mathbf{w}}^T \mathbf{x}$

4.2 Convex Function: $f(t x_1 + (1-t)x_2) \leq t f(x_1) + (1-t)f(x_2)$

4.3 Gradient Descent: 4.4 GD for Linear Regression:

```

procedure GD( $\mathcal{D}, \theta^{(0)}$ ) \mathcal{J}(\mathbf{D}) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2
1:  $\theta \leftarrow \theta^{(0)}$ 
2: while not converged do
3:   for  $i \in \{1, 2, \dots, N\}$  do
4:     for  $k \in \{1, 2, \dots, K\}$  do
5:        $\theta_k \leftarrow \theta_k - \lambda \frac{d}{d \theta_k} J^{(i)}(\theta)$ 
6:     return  $\theta$ 

```

Note: Expectation of a randomly chosen $\mathbf{V} \mathbf{J}^{(i)}(0)$ is $\mathbf{V} \mathbf{J}(0)$

Note: Steps to converge: GD O(ln 1/error), SGD O(1/error).

Note: Computation per iteration: GD O(NM), SGD O(M)

Logistic Regression

5.1 Logistic Regression and its SGD

Learning:

θ ← θ