

Algorithms

MinSeok Song

Spectral clustering

Definition 1. $Ncut(A, B)$ is defined as

$$Ncut(A, B) = cut(A, B) \left(\frac{1}{d(A)} + \frac{1}{d(B)} \right)$$

where $d(A) = \sum_{i \in A} d_i$.

Remark 1. • The intuition is that when A is relatively small, the $\frac{1}{d(A)}$ will be large, hence discouraging the isolating small groups.

- Finding minimum Ncut is equivalent to finding vector v that minimizes

$$\frac{v^T L v}{v^T D v} \text{ such that } v^T D 1 = 0, v_i \in \{a, b\}$$

where $L = D - W$.

Algorithms for Linearly Separable Data

There are several ways. First, linear programming.

Algorithm 1 Linear Programming for Classifier

Objective: minimize $\mathbf{u} = (0, \dots, 0)$ (dummy variable)

Subject to: $A\mathbf{w} \geq \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$

$A_{i,j} = y_i x_{i,j}$ (where j 'th element of the vector x_i)

Remark 2. • u is a dummy variable here; we essentially only check if the constraint is satisfied.

- This is only applied for when the data is separable.

- We can formula the regression problem with loss function $l(h, (x, y)) = |h(x) - y|$ using linear programming.

Algorithm 2 Batch Perceptron

```

1: function BATCHPERCEPTRON( $x_1, y_1, \dots, x_m, y_m$ )
2:    $w(1) \leftarrow (0, \dots, 0)$ 
3:   for  $t \leftarrow 1, 2, \dots$  do
4:     if there exists  $i$  such that  $y_i \langle w(t), x_i \rangle \leq 0$  then
5:        $w(t+1) \leftarrow w(t) + y_i x_i$ 
6:     else
7:       output  $w(t)$ 
8:       break
9:     end if
10:  end for
11: end function

```

Remark 3. • Note that $y_i(w^{(t+1)}, x_i) = y_i(w^{(t)}, x_i) + \|x_i\|^2$

- The algorithm must stop after at most $(RB)^2$ iterations, where $R = \max_i \|x_i\|$ represents a data spread, and $B = \min\{\|w\| : i \in [m], y_i \langle w, x_i \rangle \geq 1\}$ represents margin.
- To prove this, it suffices to show that $1 \geq \frac{\langle w^*, w^{(T+1)} \rangle}{\|w^*\| \|w^{(T+1)}\|} \geq \frac{\sqrt{T}}{RB}$, which we proceed by bounding numerator and denominator separately.
- We can prove that this bound is tight. For some vector $w^* \in \mathbb{R}^d$, the algorithm incurs $m = (BR)^2$ errors (considering $m = d$).
- Moreover, for $d = 3$, an algorithm can be designed to commit exactly (m) errors for any given $m \in \mathbb{N}$, serving as an upper bound concurrently

Logistic Regression

Definition 2. Fit the logistic function $\phi_{sig}(x) = \frac{1}{1 + \exp(-\langle w, x \rangle)}$ with minimization scheme $w = \operatorname{argmin}_{w \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y_i \langle w, x_i \rangle))$.

Remark 4. • The explanatory variable is between 0 and 1, making it interpretable as a probability.

- Appropriate for binary classification.
- Logistic loss function is a convex function so it's efficient to minimize.

Variational Inference and EM algorithm

Definition 3. *Kullback-Leibler(KL) divergence between p.d.f.s g and f is given by*

$$d_{KL}(g\|f) = E_g[\log(\frac{g}{f})]$$

This is always nonnegative and it can be shown by Jensen's inequality. Intuitively, we have more 'confidence' on g whenever g is greater than f , whence the logarithm is positive.

Definition 4. *The ELBO (Evidence Lower-Bound) of a p.d.f. g with respect to an unnormalized p.d.f. \tilde{f} is given by*

$$ELBO(g) := E_g[\log \frac{\tilde{f}}{g}]$$

Remark 5. • Simple algebra yields $ELBO(g) \leq \log c$, where c is a normalizing constant.

- Let's think in Bayesian framework; our unnormalized function in this case is $\tilde{f}(\theta) := f(y|\theta)f(\theta)$, so $ELBO(g) \leq \log f(y)$.
- This justifies the name "evidence lower-bound" and this helps with the choice of modeling (essentially maximizing ELBO).
- We write $ELBO(g)$, but really, what's omitted is that this is with respect to $\tilde{f}(\theta)$.

Remark 6. • It follows that $ELBO(g) = E_g[\log f(y|\theta)] - d_{KL}(g(\theta)\|f(\theta))$

- In another Bayesian framework, similar equality (will be used in EM algorithm) is

$$\begin{aligned} ELBO(g, \theta) &= \int \log\left(\frac{f(y, z|\theta)}{g(z)}\right)g(z)dz \\ &= -d_{KL}(g(z)\|f(z|y, \theta)) + \log(f(y|\theta)) \end{aligned}$$

The first term promotes matching the data, and the second term promotes matching prior beliefs. The reason we work with log domain is that, except for obvious reasons, it helps with numerical stability. Indeed, remark 2 motivates what's called EM algorithm.

Theorem 1. *We have $\log f(y|\theta_l) \leq \log f(y|\theta_{l+1})$*

Proof. Let $g_l(z) = f(z|y, \theta_l)$. We have

Algorithm 3 EM Algorithm

- 1: **Input:** Initialization of θ_0
- 2: **repeat**
- 3: **E-step:** compute

$$E_{Z \sim f(z|y, \theta_l)}[\log f(y, Z | \theta)] = \int \log f(y, z | \theta) f(z | y, \theta_l) dz$$

- 4: **M-step:** compute

$$\theta_{l+1} = \arg \max_{\theta} E_{Z \sim f(z|y, \theta_l)}[\log f(y, Z | \theta)]$$

- 5: **until** some stopping criterion
-

$$\begin{aligned} \log f(y|\theta_{l+1}) &= ELBO(g_l, \theta_{l+1}) + d_{KL}(g_l | f(z|y, \theta_{l+1})) \\ &\geq ELBO(g_l, \theta_l) + d_{KL}(g_l | f(z|y, \theta_l)) \\ &= \log f(y|\theta_l) \end{aligned}$$

□

- This shows that likelihood function is non-decreasing for each iteration, and since likelihood is always bounded by 1, we have established the convergence of the algorithm.
- GMM algorithm is a specific instance of this algorithm. Here, w_{ik} can be thought of as latent variable, corresponding to E-step, and computing θ and π corresponds to M-step.
- E-step can be computationally expensive and so we usually use Monte-Carlo method to approximate.

Remark 7. • Let us now venture into variational inference, which we use when approximating the intractable distribution. For the choice of possible sets that f admits in $d_{KL}(g \| f)$, we can use **mean field family**, which assumes the independence for coordinate distributions.

Theorem 2. Let $g(x) = \prod_{i=1}^d g_i(x_i)$ with $g_{-i}(x_{-i})$ fixed. Then

$$g_i^*(x_i) \propto \exp(E_{g_{-i}}[\log f(x_i, x_{-i})])$$

maximizes $ELBO(g)$.

Proof. By routine algebra (we need to use independence at some point), one can show that

$$ELBO(g) = E_{g_i}[\log(\exp(E_{g_{-i}}[\log f(x_i, x_{-i})])) - \log g_i(x_i)] + C$$

and notice that the first term can be phrased as $-d_{KL}(g^* \| g)$, whence $g^* = g$ gives the optimization solution. □

- Remark 8.*
- The intuition is to average out the effect of x_{-i} on log expected value in order to incorporate the independence between g_i 's.
 - This leads to the CAVI algorithm, which approximates the unnormalized target density \tilde{f} . After initialization, updating g_i will increase ELBO for each i , so may iterate until ELBO converges.
 - The disadvantage is that it may be computationally expensive and accuracy might be not so good.