

10/19 Numerical optimization

How to minimize (maximize) a function $f(\vec{\theta})$ over $\vec{\theta} = (\theta_1, \dots, \theta_d) \in \mathbb{R}^d$?

1) Gradient descent (or ascent)

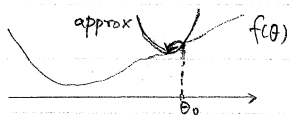
$$\vec{\theta}' \leftarrow \vec{\theta} - \eta \left(\frac{\partial f}{\partial \vec{\theta}} \right)$$

scalar learning rate $\eta > 0$.

2) Newton's method

* Approximate $f(\theta)$ near $\theta = \theta_0$.

$$f(\theta) \approx f(\theta_0) + f'(\theta_0)(\theta - \theta_0) + \frac{1}{2} f''(\theta_0)(\theta - \theta_0)^2 + \dots$$



* Minimize quadratic approximation.

$$\frac{\partial}{\partial \theta} \left[f(\theta_0) + f'(\theta_0)(\theta - \theta_0) + \frac{1}{2} f''(\theta_0)(\theta - \theta_0)^2 \right] = 0$$

$$f'(\theta_0) + f''(\theta_0)(\theta - \theta_0) = 0$$

$$\theta^* = \theta_0 - \frac{f'(\theta_0)}{f''(\theta_0)}$$

* Update rule

$$\theta \leftarrow \theta - \frac{f'(\theta)}{f''(\theta)}$$

* In d dimensions:

Define $H = d \times d$ matrix of 2nd partial derivatives

$H_{ij} = \frac{\partial^2 f}{\partial \theta_i \partial \theta_j}$ (symmetric) Hessian matrix

Update rule: $\vec{\theta} \leftarrow \vec{\theta} - H^{-1} \left(\frac{\partial f}{\partial \vec{\theta}} \right)$ evaluated at current value of $\vec{\theta}$

matrix inverse of Hessian

matrix-vector multiplication.

* Pros

- no learning rate that needs to be tuned.
- converges very fast (when it converges)

* Cons

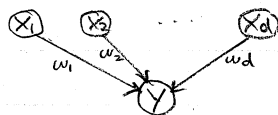
- unstable if far from optimum
- often expensive to compute or invert Hessian
 $O(d^2)$ $O(d^3)$
- converges only to local minimum (if it converges).

Case IIb. Learning in BNs with sigmoid CPTs, complete data (logistic regression)

* belief network

parents $\vec{X} \in \mathbb{R}^d$

child $Y \in \{0, 1\}$



* Sigmoid CPT

$$P(Y=1 | \vec{X}=\vec{x}) = \sigma(\vec{w} \cdot \vec{x}) \quad \text{with} \quad \sigma(z) = \frac{1}{1 + \exp(-z)}$$

* properties of sigmoid function

$$\sigma(-z) = 1 - \sigma(z)$$

$$\sigma'(z) = \sigma(z) \sigma(-z)$$

* log-likelihood of training examples

$$P(\vec{X}_t, y_t) \Big|_{t=1}^T \quad \text{with} \quad y_t \in \{0, 1\}$$

$$\mathcal{L}(\vec{w}) = \log P(\text{data}) \quad \text{assume data is i.i.d.}$$

$$= \log \prod_{t=1}^T P(Y=y_t | \vec{X}=\vec{x}_t)$$

$$= \sum_{t=1}^T \log P(Y=y_t | \vec{X}=\vec{x}_t)$$

$$= \sum_{t=1}^T \log [\sigma(\vec{w} \cdot \vec{x}_t)^{y_t} \sigma(-\vec{w} \cdot \vec{x}_t)^{1-y_t}]$$

$$= \sum_{t=1}^T [y_t \log \sigma(\vec{w} \cdot \vec{x}_t) + (1-y_t) \log \sigma(-\vec{w} \cdot \vec{x}_t)]$$

* To maximize $\mathcal{L}(\vec{w})$

$$\begin{aligned} 0 = \frac{\partial \mathcal{L}}{\partial w_\alpha} &= \sum_t \left[y_t \frac{1}{\sigma(\vec{w} \cdot \vec{x}_t)} \sigma(\vec{w} \cdot \vec{x}_t) \sigma(-\vec{w} \cdot \vec{x}_t) X_{\alpha t} \right. \\ &\quad \left. + (1-y_t) \frac{1}{\sigma(-\vec{w} \cdot \vec{x}_t)} \sigma(-\vec{w} \cdot \vec{x}_t) \sigma(\vec{w} \cdot \vec{x}_t) (-X_{\alpha t}) \right] \\ &= \sum_t [y_t (1 - \sigma(\vec{w} \cdot \vec{x}_t)) X_{\alpha t} - (1-y_t) \sigma(\vec{w} \cdot \vec{x}_t) X_{\alpha t}] \\ &= \sum_t [y_t - \sigma(\vec{w} \cdot \vec{x}_t)] X_{\alpha t} \quad \text{for } \alpha = 1, 2, \dots, d \end{aligned}$$

* Gradient

$$\frac{\partial \mathcal{L}}{\partial \vec{w}} = \sum_t (y_t - \sigma(\vec{w} \cdot \vec{x}_t)) \vec{x}_t$$

difference between target value $y_t \in \{0, 1\}$ and model's prediction $P(Y=1 | \vec{x}_t)$

* Hessian

$$H_{\alpha\beta} = \frac{\partial^2 \mathcal{L}}{\partial w_\alpha \partial w_\beta} = - \sum_t \sigma(\vec{w} \cdot \vec{x}_t) \sigma(-\vec{w} \cdot \vec{x}_t) X_{\alpha t} X_{\beta t}$$

$$H = \frac{\partial^2 \mathcal{L}}{\partial \vec{w} \partial \vec{w}^T} = - \sum_t \sigma(\vec{w} \cdot \vec{x}_t) \sigma(-\vec{w} \cdot \vec{x}_t) \vec{x}_t \vec{x}_t^T$$

* Algorithms for ML estimation.

1) Gradient ascent : update $\vec{w} \leftarrow \vec{w} + \eta \left(\frac{\partial \mathcal{L}}{\partial \vec{w}} \right)$ suggest: $\eta = \frac{0.02}{T}$

2) Newton's method : update $\vec{w} \leftarrow \vec{w} - H^{-1} \left(\frac{\partial \mathcal{L}}{\partial \vec{w}} \right)$

* global optimality

It can be shown that $\mathcal{L}(\vec{w})$ for logistic regression has no spurious local maxima. i.e., $\mathcal{L}(\vec{w})$ is convex.

Case III fixed DAG, discrete nodes, "lookup" CPTs, incomplete data

* Variables in BN : H = hidden variables } may vary from one example to next.

V = visible variables

$$X = HUV$$

* Data set: Assume T incomplete/partial examples drawn i.i.d from P(X).

#	X ₁	X ₂	...	X _n
1	5	9		?
2	?	4		2
3	?	?		1
⋮				
T	1	2		2

* log-likelihood

$$\begin{aligned}
 \mathcal{L} &= \log P(\text{DATA}) && \text{assume i.i.d.} \\
 &= \log \prod_{t=1}^T P(V^{(t)}) && \text{observed nodes in BN for } t^{\text{th}} \text{ example.} \\
 &= \sum_T \log P(V^{(t)}) \\
 &= \sum_T \log \sum_{h^{(t)}} P(v^{(t)}, h^{(t)}) && \text{joint distribution} \\
 &= \sum_T \log \sum_{h^{(t)}} \prod_{i=1}^n P(X_i | pa(X_i)) \Big|_{\substack{V=v^{(t)} \\ H=h^{(t)}}}
 \end{aligned}$$

Before for complete data:

CPTs decoupled ⇒ many independent optimizations

Now: all CPTs are coupled.

How to optimize?

* Expectation - Maximization (EM)

- alternative to gradient ascent or Newton's method.

* Intuition - by analogy

• ML estimates for complete data

$$I(a,b) = \begin{cases} 1 & \text{if } a=b \\ 0 & \text{if } a \neq b \end{cases} \text{ indicator function.}$$

$$P_{ML}(X_i = x | pa(X_i) = \pi) = \frac{\text{count}(X_i = x, pa_i = \pi)}{\text{count}(pa_i = \pi)} = \frac{\sum_T I(X_i^{(t)}, x) I(pa_i^{(t)}, \pi)}{\sum_T I(pa_i^{(t)}, \pi)}$$

• For incomplete data, we must "fill in" missing values:

$$P(X_i = x | pa_i = \pi) \leftarrow \frac{\sum_T P(X_i = x, pa_i = \pi | V^{(t)})}{\sum_T P(pa_i = \pi | V^{(t)})}$$

RHS reduces to previous formula for complete data.

Intuition: expected statistics under P(H|V) substitute for observed statistics in complete data case.