# @ Lecture 10

Case IIB   Sigmoid CPTs ( logistic regression)

$$\vec{X} \in \mathbb{R}^d \quad \textcircled{X_1} \quad \textcircled{X_2} \quad \cdots \quad \textcircled{X_d}$$



$$Y \in \{0,1\} \quad \textcircled{Y}$$

- sigmoid CPT
$$P(Y=1|\vec{X}) = \sigma(\vec{\omega} \cdot \vec{X})$$
   with $\sigma(z) = \dfrac{1}{1+e^{-z}}$

properties of $\sigma(z)$:

- $\sigma(z) = 1 - \sigma(-z)$
$$\frac{d}{dz}\sigma(z) = \sigma(z)\,\sigma(-z)$$

\* **Training Examples** (IID)
$$\{(\vec{X_t}, Y_t)\}_{t=1}^{T}$$

- log (conditional) likelihood :
$$\mathcal{L}(\vec{\omega}) = \log P(data)$$

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

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

$$\mathcal{L}(\vec{w}) = \sum_{t=1}^{T} \left[ \sigma(\vec{w} \cdot \vec{x}_t) \right]^{y_t}$$

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

with $y_t \in \{0,1\}$

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

To maximize this expression:

$$0 = \frac{\partial \mathcal{L}}{\partial w_\alpha} = \sum_t \left[ y_t \frac{1}{\sigma(\vec{w} \cdot \vec{x}_t)} \sigma(\vec{w} \cdot \vec{x}_t) \sigma(-\vec{w} \cdot \vec{x}_t) x_{\alpha t} + \right.$$
$$\left. (1-y_t) \frac{1}{\sigma(-\vec{w} \cdot \vec{x}_t)} \sigma(\vec{w} \cdot \vec{x}_t) \sigma(-\vec{w} \cdot \vec{x}_t)(-x_{\alpha t}) \right]$$

$$= \sum_t x_{\alpha t} \left[ y_t \, \sigma(-\vec{w} \cdot \vec{x}_t) - (1-y_t) \sigma(\vec{w} \cdot \vec{x}_t) \right]$$

$$= \sum_t x_{\alpha t} \left[ y_t (1 - \sigma(\vec{w} \cdot \vec{x}_t)) - (1-y_t) \sigma(\vec{w} \cdot \vec{x}_t) \right]$$

$$= \sum_t x_{\alpha t} \left[ y_t - \sigma(\vec{w} \cdot \vec{x}_t) \right]$$

difference$^{\underline{i}}$ between target value $y \in \{0,1\}$
and $P(Y=1 | \vec{x}_t)$ modeling
our prediction.

$$0 = \frac{\partial \mathcal{L}}{\partial w_\alpha} \quad \text{for} \quad \alpha = 1, 2, \ldots, d$$

These are NON-LINEAR equations.

- Hessian Matrix:

$$H_{\alpha\beta} = \frac{\partial^2 \mathcal{L}}{\partial w_\alpha w_\beta} = \frac{\partial}{\partial w_\beta} \left\{ \sum_t \left[ y_t - \sigma(\vec{w} \cdot \vec{x}_t) \right] x_{\alpha t} \right\}$$

$$= -\sum_t \sigma(\vec{w} \cdot \vec{x}_t) \, \sigma(-\vec{w} \cdot \vec{x}_t) \, x_{\beta t} \, x_{\alpha t}$$

Vector form:

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

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

ML Estimation:

1) Gradient ~~descent~~ Ascent : update $\vec{w} \leftarrow \vec{w} + \eta \frac{\partial \mathcal{L}}{\partial \vec{w}}$

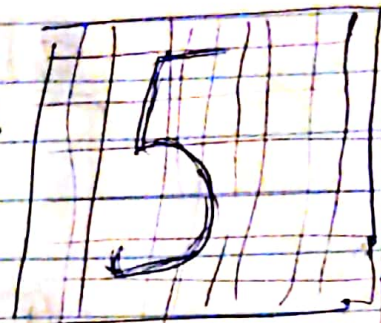$$\text{suggest:} \quad \eta = \frac{0.2}{T}$$

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

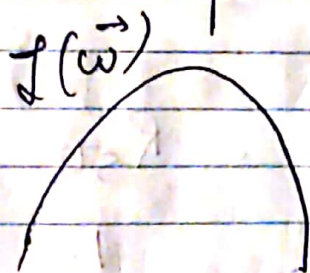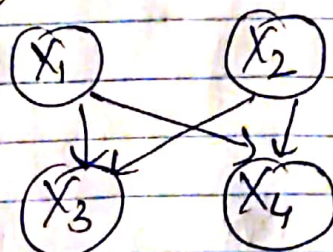$$\text{suggest} \quad \vec{w} = [0, \ldots 0]^T$$

HW5: Classify

$\boxed{3}$ vs $\boxed{5}$

                y=0                     y=1

- Global Optimality: ← it can be shown that $\mathcal{L}(\vec{\omega})$ is concave for logistic regression. and has no spurious local maxima.

$\mathcal{L}(\vec{\omega})$

Case III: fixed DAG, discrete nodes, lookup CPT's,

INCOMPLETE DATA.

TOY Example:

$X_1$       $X_2$     binary nodes $X_i \in \{0, 1\}$

$X_3$       $X_4$

| t | $X_1$ | $X_2$ | $X_3$ | $X_4$ |
|---|---|---|---|---|
| 1 | ? | 1 | 0 | 1 |
| 2 | 0 | ? | 1 | 1 |
| 3 | ? | ? | 0 | 1 |
| ⋮ | | | | |
| T | | | | |

## Ex: Movie Recommender System:

$Z \in \{1, 2, \ldots K\}$ types of movie you?

$R_i \in \{0, 1\}$ movie rating

i = 1  Avengers
i = 2  Toy Story
i = 3  Star Wars ... n.

| Z | $R_1$ | $R_2$ | | | $R_{50}$ |
|---|---|---|---|---|---|
| 1 | ? | 1 | ? | | 1 |
| 2 | ? | ? | 0 | | 0 |
| 3 | ? | ? | ? | 1 | 1 |
| ⋮ | | | | | |
| 256 | ? | | | | |

↓
no. of students

* Variables in BN

$H$ = set of Hidden (unobserved) variables
$V$ = set of visible (observed) variables

Can vary with example

$X = H \cup V$ (all nodes)

* **log-likelihood:**

$$\mathcal{L} = \log P(\text{data})$$

$$= \log \prod_{t=1}^{T} P(V^{(t)}) \quad \overset{\text{marginal}}{\underset{\text{probs.}}{}}$$

$$\hookrightarrow \text{visible nodes on } t^{th} \text{ example}$$

$$= \sum_t \log P(V^{(t)})$$

$$= \sum_t \log \sum_h P(H=h, V^{(t)}) \quad \text{marginalization.}$$

$$\mathcal{L} = \sum_t \log \sum_h \left\{ \prod_{i=1}^{n} P(X_i = x \mid pa_i = \pi) \right\} \Bigg|_{\substack{H=h \\ V=V^{(t)}}} \overset{\text{in form}}{\underset{\text{of CPTs}}{}}$$

For complete data:
    CPTs decoupled $\Leftrightarrow$ many independent optimizations.

Now, for incomplete data:
    many (or all) CPTs are coupled.

How to optimize?

Options?

1) Gradient Descent: $\vec{\theta} \leftarrow \vec{\theta} + \eta \dfrac{dL}{d\vec{\theta}}$ must Tune $\eta > 0$ asymptotic but not monotonic convg.

2) Newton's Method: $\vec{\theta} \leftarrow \vec{\theta} - H^{-1} \dfrac{dL}{d\vec{\theta}}$ expensive, fast but unstable!

3) New method: Auxilliary functions $Q(\vec{\theta}, \vec{\theta}')$.
   How to minimize $f(\vec{\theta})$?

- Suppose $Q(\vec{\theta}, \vec{\theta}')$ satisfies two properties

EQUALITY $\quad (1.) \; Q(\vec{\theta}, \vec{\theta}) = f(\vec{\theta})$

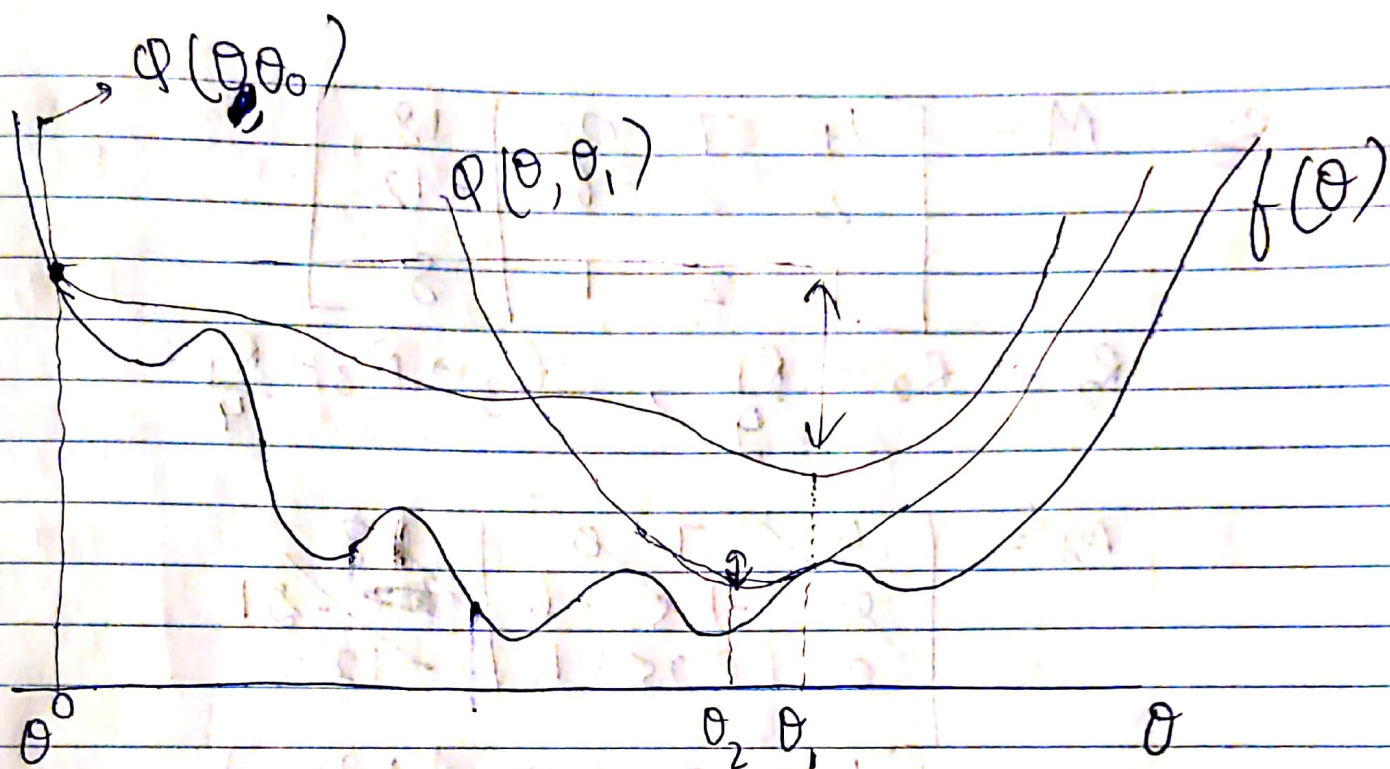BOUND $\quad (11.) \; Q(\vec{\theta}', \vec{\theta}) \geq f(\vec{\theta}') \; \forall \vec{\theta}, \vec{\theta}'$

Consider update rule:

$$\vec{\theta}_{new} = \underset{\vec{\theta}}{\arg\min} \; Q(\vec{\theta}, \vec{\theta}_{old})$$

Now, $f(\vec{\theta}_{new}) \leq Q(\vec{\theta}_{new}, \vec{\theta}_{old})$ by property $(11.)$

$\qquad\qquad \leq Q(\vec{\theta}_{old}, \vec{\theta}_{old})$ by update rule

$\qquad\qquad = f(\vec{\theta}_{old})$ by property $(1.)$

By iterating: $f(\vec{\theta}_0) \geq f(\vec{\theta}_1) \ldots \ldots \geq f(\vec{\theta}_n)$

$\varphi(\theta, \theta_0)$

$\varphi(\theta, \theta_1)$

$f(\theta)$

$\theta^0$   $\theta_2 \, \theta_1$   $\theta$

- Properties:

- no learning rate:
- monotonic improvement.
- convergence to stationary point where
  gradient vanishes. $\frac{\partial f}{\partial \theta}$