

Logistic regression:

①

- Nice, useful model binary classifier

Model

$$\log \frac{P(1|x)}{P(-1|x)} = w^T x$$

↑

you can have a constant here by extending x

So

$$P(1|x) = \frac{e^{w^T x}}{1 + e^{w^T x}}$$

We can fit with maximum likelihood

$$L(w) = \sum_{i \in \text{examples}} \log P(y_i | x_i)$$

this is equiv to minimizing -ve log-likelihood ⁽²⁾

$$w = \underset{w}{\operatorname{argmin}} - \sum_{i \in \text{examples}} \log P(y_i | x_i)$$

$$= \underset{w}{\operatorname{argmin}} - \sum_i \left[\frac{(y_i + 1)}{2} \cdot w^T x_i - \log(1 + e^{w^T x_i}) \right]$$

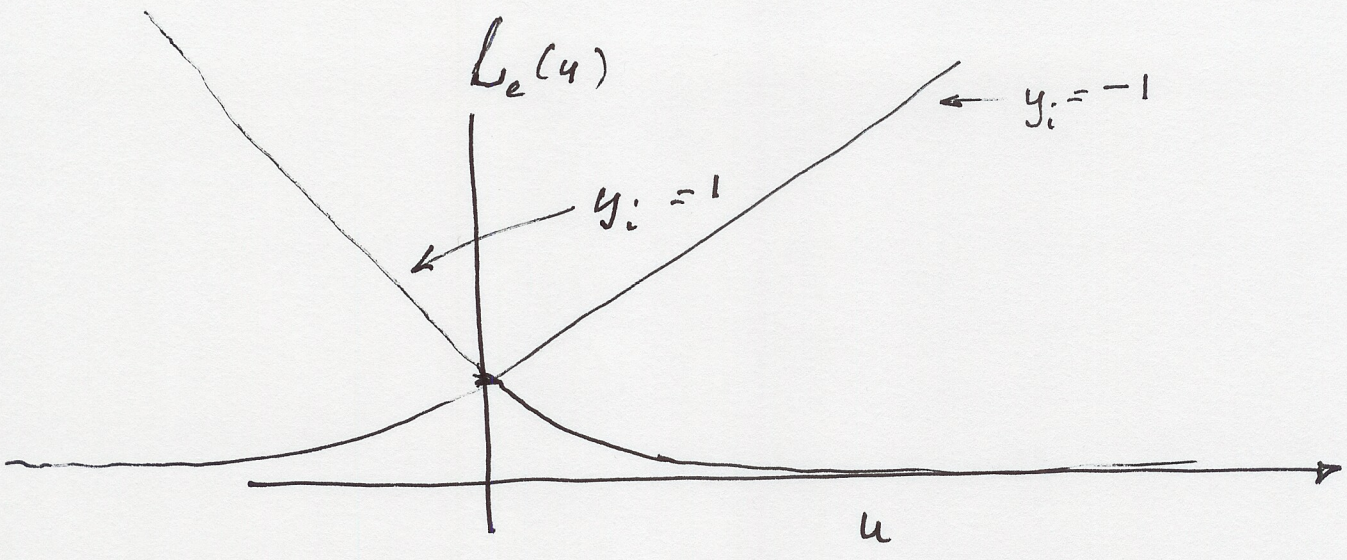
Now regard

$$\log(1 + e^{u_i}) - \frac{(y_i + 1)}{2} u_i \text{ as a loss}$$

$$L(y_i, u_i) \quad \leftarrow \text{exponential loss}$$

↑ ↑

Value of example score = $w^T x_i$



looks a lot like hinge loss :

So we can think of Lik. as minimizing a loss function.

How? $\nabla = 0$, Newton.

$$\nabla_w L_e = \sum_i \left[\frac{x_i}{1 + e^{w^T x_i}} - \frac{(y_i + 1) x_i}{2} \right]$$

$$H_w L_e = \sum_i \frac{e^{w^T x_i}}{(1 + e^{w^T x_i})^2} \cdot x_i x_i^T$$

④

Notice:

- examples where $w^T x_i$ has large abs value have little effect on H .
- For ~~these~~ others, H looks like a covariance
- what if features are correlated?
 - H has small eigenvalues
 - ests of w will be unreliable in these dirs.
 - Should manifest as large w

⇒ Regularize.

Solve

argmin
 ω

$L_e(\omega, \text{examples}) + \lambda \|\omega\|$

?

Some work