

Classification

probabilistic classifiers

→ logistic regression

We consider a dataset

$$\{x^{(i)}, t^{(i)}\}_{i=1}^N$$

$$t^{(i)} \in \{0, 2\}$$

$$h_{\beta}(x) = \sigma(\beta_0 + \beta_1 x_1 + \dots + \beta_D x_D)$$

if we take $\sigma(x) = \frac{1}{1 + e^{-x}}$ (sigmoid)

$\Rightarrow h_{\beta}(x) \in [0, 1]$ and can be interpreted as a probability

let us define $p(t(x)=1 | x) = h_{\beta}(x)$

$$p(t(x)=0 | x) = 1 - h_{\beta}(x)$$

to train our logistic regression classifier,
let us consider

$$p(t(x) = t^{(i)} | x) = \underbrace{h_{\beta}(x)^{t^{(i)}} (1 - h_{\beta}(x))^{1 - t^{(i)}}}$$

let us consider the whole dataset

$$p(\underbrace{\{t(x^{(i)}) = t^{(i)}\}_{i=1}^N}_{\text{dataset}} | \{x^{(i)}\}_{i=1}^N) = \prod_{i=1}^N \underbrace{h_{\beta}(x^{(i)})^{t^{(i)}} (1 - h_{\beta}(x^{(i)}))^{1 - t^{(i)}}}_{\text{likelihood}}$$

$$\log \left(p(\{t(x^{(i)}) = t^{(i)}\}_{i=1}^N | x^{(i)}) \right)$$

$$= \sum_{i=1}^N t^{(i)} \log h_{\beta}(x^{(i)}) + (1 - t^{(i)}) \log(1 - h_{\beta}(x^{(i)}))$$

We can then train the model by minimizing $-\log$,

$$\beta^* = \operatorname{arg\,min}_{\beta} \underbrace{\sum_{i=1}^N t^{(i)} \log h_{\beta}(x^{(i)}) + (1 - t^{(i)}) \log(1 - h_{\beta}(x^{(i)}))}_{\parallel}$$

→ can be solved through gradient descent $\mathcal{L}(\beta)$

for any β_j , for any sample pair $\{t^{(i)}, x^{(i)}\}$

$$\begin{aligned} \frac{\partial \ell_i}{\partial \beta_j} &= \overbrace{t^{(i)}} \frac{\partial}{\partial \beta_j} \log(\sigma(\beta^T \tilde{x}^{(i)})) + (1-t^{(i)}) \frac{\partial}{\partial \beta_j} \log(1-\sigma); \\ &= \overbrace{t^{(i)}} \left(\frac{\frac{\partial}{\partial \beta_j} \sigma(\beta^T \tilde{x}^{(i)})}{\sigma(\beta^T \tilde{x}^{(i)})} \right) + (1-t^{(i)}) \left(- \frac{\frac{\partial}{\partial \beta_j} \sigma(\beta^T \tilde{x}^{(i)})}{(1-\sigma(\dots))} \right) \end{aligned}$$

using the chain rule

$$\frac{\partial}{\partial \beta_j} \sigma(\beta^T \tilde{x}^{(i)}) = \underbrace{\sigma'(\beta^T \tilde{x}^{(i)})}_{\text{green circle}} \cdot \underbrace{\tilde{x}_j^{(i)}}_{\text{blue circle}}$$

$$\frac{d\sigma(x)}{dx} = \frac{d}{dx} \frac{1}{1+e^{-x}} = \frac{e^{-x}}{(1+e^{-x})^2} = \left(1 - \frac{1}{1+e^{-x}}\right) \frac{1}{1+e^{-x}}$$

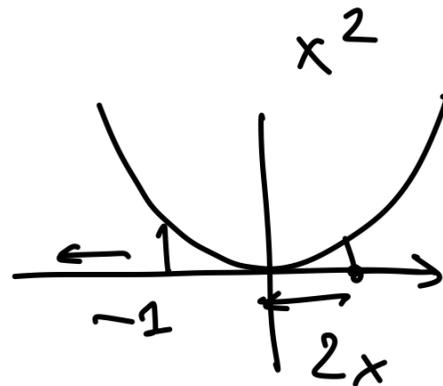
$$\sigma'(x) = \sigma(x)(1-\sigma(x))$$

$$\frac{\partial \sigma(\beta^T \tilde{x}^{(i)})}{\partial \beta_j} = \tilde{x}_j^{(i)} \sigma(\beta^T \tilde{x}^{(i)}) (1 - \sigma(\beta^T \tilde{x}^{(i)}))$$

$$\frac{\partial \lambda^{(i)}}{\partial \beta_j} = t^{(i)} \tilde{x}_j^{(i)} \frac{\sigma(1-\sigma)}{\sigma} - (1-t^{(i)}) \cdot \frac{1}{1-\sigma} \tilde{x}_j^{(i)} \sigma(1-\sigma)$$

$$= t^{(i)} \tilde{x}_j^{(i)} (1-\sigma) - (1-t^{(i)}) \sigma \tilde{x}_j^{(i)}$$

$$= (t^{(i)} - \sigma) \tilde{x}_j^{(i)}$$



Now we can use update of the form

$$2x = -2$$

$$\beta_j^{(k+1)} \leftarrow \beta_j^{(k)} - \eta (t^{(i)} - \sigma(\beta^T \tilde{x}^{(i)})) \tilde{x}_j^{(i)}$$

SGD (one sample)

$$\widehat{\beta}_j^{(k+1)} \leftarrow \widehat{\beta}_j^{(k)} - \eta \sum_{i=1}^N \overbrace{(t^{(i)} - \sigma(\beta^T \tilde{x}^{(i)}))} \widehat{\tilde{x}}_j^{(i)}$$

BATCH GD

$$\tilde{x} = [1, \vec{x}] \rightarrow \in \mathbb{R}^{D+1}$$

$$x \in \mathbb{R}^D$$

Classification

→ least squares

→ logistic regression

→ Gaussian Discriminant Analysis

→ perceptron

general formulation $\sigma(\beta_0 + \beta_1 x_1 + \dots + \beta_D x_D)$

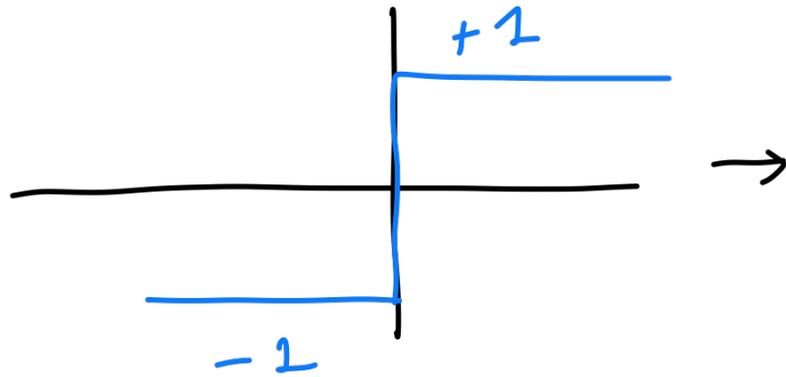
least squares: $\sigma(x) = x$

logistic regression $\sigma(x) = \frac{1}{1 + e^{-x}}$

perceptron

$$\sigma(x) = \begin{cases} +1 & x > 0 \\ -1 & x \leq 0 \end{cases}$$

Perceptron



→ derivative is zero almost everywhere

→ Solution: take loss

$$t^{(i)} = \begin{cases} +1 & x^{(i)} \in C_0 \\ -1 & x^{(i)} \in C_1 \end{cases}$$

$$l(\beta) = - \sum_{i \in \text{Misclassified}} \widehat{t^{(i)}} \cdot (\beta^T \tilde{x}^{(i)})$$

$$\min l(\beta) = \min - \sum_{i \in \text{Misclassified}} t^{(i)} (\beta^T \tilde{x}^{(i)})$$

From this we get

$$\text{grad}_{\beta} l = - \left(\sum_{i \in H} t^{(i)} \tilde{x}^{(i)} \right)$$

Perceptron update (stochastic)

While \exists misclassified points s.t. $\sigma(\beta^T \tilde{x}^{(i)}) \neq t^{(i)}$

take a misclassified $x^{(i)}$ and do

$$\beta^{(k+1)} \leftarrow \beta^{(k)} + \eta \underbrace{t^{(i)}} \underbrace{\tilde{x}^{(i)}}$$

perceptron learning rule

Perceptron convergence theorem: if dataset is linearly separable perceptron will converge in a finite number of steps. \rightarrow and η small enough

Assume $\beta_0 = 0$

