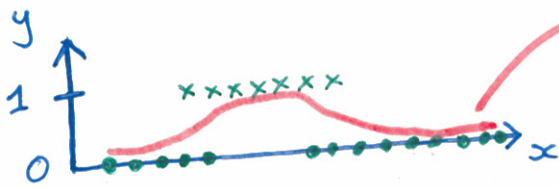


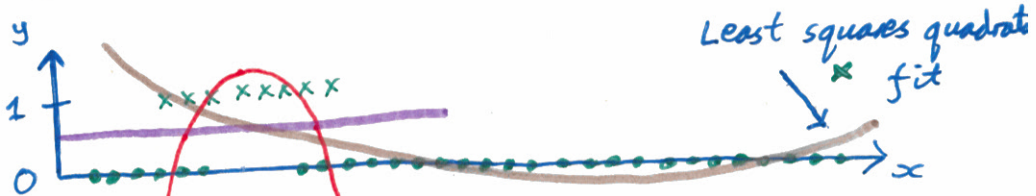
Regressing on Labels



$$f(x) \approx P(y=1|x)$$

(If enough data and basis functions, linear least squares works!)

Often bad idea:



Function can give good labels
Terrible square error

Gradients for least squares

$N \times D$ design matrix
of all inputs

$$\text{Residuals } \underline{r} = \underline{y} - \underline{X}\underline{w}$$

$$\text{Cost } \underline{r}^T \underline{r} = (\underline{y} - \underline{X}\underline{w})^T (\underline{y} - \underline{X}\underline{w})$$

$$= \underline{y}^T \underline{y} - 2\underbrace{\underline{w}^T (\underline{X}^T \underline{y})}$$

$$+ \underline{w}^T \underline{X}^T \underline{X} \underline{w}$$

"Gradient"

vector of partial derivatives

$$\nabla_{\underline{w}} [\underline{r}^T \underline{r}] = -2 \underline{X}^T \underline{y} + 2 \underline{X}^T \underline{X} \underline{w}$$

Gradient descent:

$$\underline{w} \leftarrow \underline{w} - \alpha \nabla_{\underline{w}} [\underline{r}^T \underline{r}]$$

α
step-size, small.

$$\nabla_{\underline{w}} \underline{w}^T \underline{b} = \begin{bmatrix} \frac{\partial \underline{w}^T \underline{b}}{\partial w_1} \\ \frac{\partial \underline{w}^T \underline{b}}{\partial w_2} \\ \vdots \\ \frac{\partial \underline{w}^T \underline{b}}{\partial w_n} \end{bmatrix} = \underline{b}$$

$$\frac{\partial}{\partial w_i} \sum_j w_j b_j = \frac{\partial}{\partial w_i} [b_1 w_1 + b_2 w_2 + \dots + b_i w_i + \dots + b_n w_n] = b_i$$

$$\nabla_{\underline{w}} [\underline{r}^T \underline{r}] = \underline{0}$$

$$(\underline{x}^T \underline{x}) \underline{w} = \underline{x}^T \underline{y} \quad \text{at optimal } \underline{w}$$

$$\underline{w} = \underbrace{(\underline{x}^T \underline{x})^{-1}}_{\text{Pseudo Inverse}} \underline{x}^T \underline{y}$$

Matlab:

$$(\underline{x}' * \underline{x}) \setminus (\underline{x}' * \underline{y})$$

$$\underline{x}^T \underline{x} \quad D \times D$$

$$\begin{array}{l} \underline{x}^{-1} \underline{x}^{-T} \underline{x}^T \underline{y} \\ \underline{w} = \underline{x}^{-1} \underline{y} \\ \underline{w} = \underline{x} \setminus \underline{y} \end{array} \quad \begin{array}{l} \nearrow \text{I} \\ \searrow \end{array}$$

~~1/2~~ $X^T X$ $D \times D$ Matrix

If rank $< D$ not invertible.

R puzzle

One hot encoding

Turns integer $k \rightarrow [00 \dots 0100 \dots 0]$

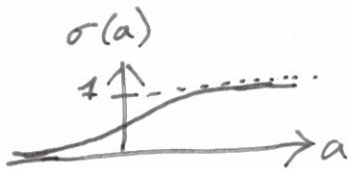
↓
in $\{1, 2, \dots, k\}$

$k-1$ long

↑
kth position

Logistic Regression

$$f(\underline{x}; \underline{w}) = \sigma(\underline{w}^T \underline{x})$$
$$= \frac{1}{1 + e^{-\underline{w}^T \underline{x}}}$$



Loss Function

Could use square loss again:

$$\sum_{n=1}^N (y^{(n)} - f(\underline{x}; \underline{w}))^2$$

Normally say

$$P(y=1 | \underline{x}) = f(\underline{x}; \underline{w})$$

and use maximum likelihood. Maximize
prob of data

Or minimize negative log probability.

$$NLL = - \sum_{n: y^{(n)}=1} \log \sigma(\underline{w}^T \underline{x}^{(n)}) - \sum_{n: y^{(n)}=0} \log(1 - \sigma(\underline{w}^T \underline{x}^{(n)}))$$

What I do is make labels $\{-1, +1\}$

$$z^{(n)} = 2y^{(n)} - 1$$

Useful $(1 - \sigma(a)) = \sigma(-a)$

$$NLL = - \sum_{n=1}^N \log \sigma(\underline{z} \underline{w}^T \underline{x}^{(n)})$$

Prob of being correct, σ_n

$$\nabla_{\underline{w}} NLL = - \sum_{n=1}^N \nabla_{\underline{w}} \log \sigma_n$$

$$= - \sum_{n=1}^N \frac{1}{\sigma_n} \nabla_{\underline{w}} \sigma_n$$

$$= - \sum_{n=1}^N \frac{1}{\sigma_n} \sigma_n (1 - \sigma_n) \nabla_{\underline{w}} z^{(n)T} \underline{x}^{(n)}$$

$$= - \sum_{n=1}^N (1 - \sigma_n) z^{(n)T} \underline{x}^{(n)}$$

