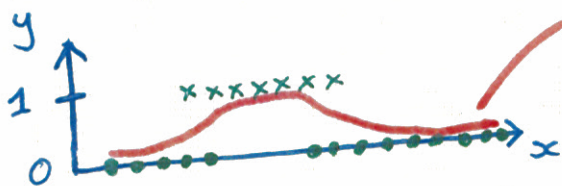


Regressing on Labels



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

(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

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

X is $N \times D$ design matrix of inputs
 \underline{w} is $D \times 1$ parameters
 \underline{y} is $N \times 1$ training labels

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

$$= \underline{y}^T \underline{y} - 2\underline{w}^T X^T \underline{y} + \underline{w}^T X^T X \underline{w}$$

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

$$\left[= -2X^T \underbrace{(\underline{y} - X\underline{w})}_{\underline{r}} \right]$$

Gradient descent:

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

η is step-size, "small" 0.01?

Normal Equations approach

$\nabla_{\underline{w}} [\underline{r}^T \underline{r}] = \underline{0}$ at least min. ~~sq~~ squares solution

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

So if the best \underline{w} is unique:

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

$$\left| \begin{array}{l} X^{-1} X^{-T} X^T y \\ \underline{w} = X^{-1} y \end{array} \right. \quad X$$

" $\underline{w} = X \setminus y$ "

Matlab:

~~$$\text{inv}(X' * X) * (X' * y)$$~~

$$(X' * X) \setminus (X' * y)$$

More accurate than

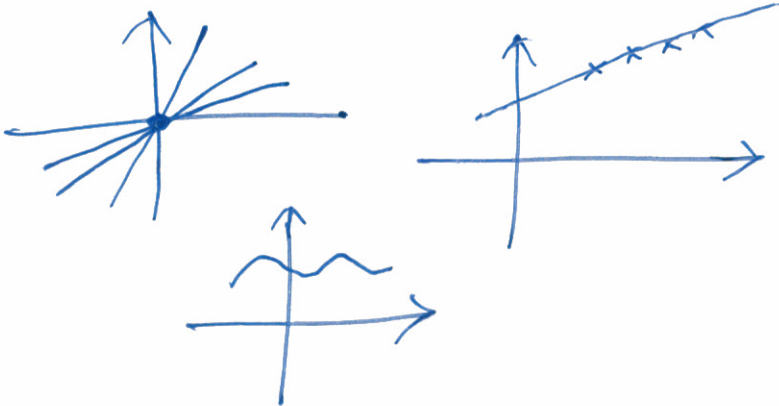
R puzzle

"red"	→	1 0 0	why?	bias
"blue"	→	0 1 0		1
"green"	→	0 0 1		1

For unique \underline{w} to min. sum of squares

$(X^T X)$ is invertible

→ Full rank.

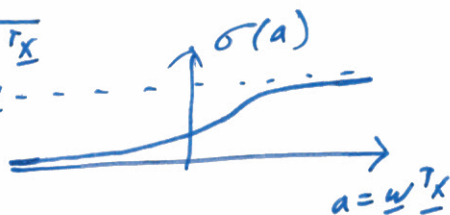


Logistic Regression

could $\phi(\underline{x})$

$$f(\underline{x}; \underline{w}) = \sigma(\underline{w}^T \underline{x}) \quad f \in [0, 1]$$

$$= \frac{1}{1 + e^{-\underline{w}^T \underline{x}}}$$



Loss Function

Could use square loss again:

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

Normal interpretation:

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

Could maximize likelihood of the parameters.

Likelihood prob. of data given the parameters.

$$\text{Prob } P(\{y^{(n)}\} | \{X\}, \underline{w}) = \prod_1 p(y^{(n)} | \underline{x}^{(n)}, \underline{w})$$

Or minimize negative log probability

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

I like to make labels $\{-1, +1\}$

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

Useful fact:

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

$$NLL = - \sum_{n=1}^N \log \underbrace{\sigma(z^{(n)} \underline{w}^T \underline{x}^{(n)})}_{\text{Prob. of being correct, } \sigma_n}$$

Prob. of being correct, σ_n

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

$$= - \sum_{n=1}^N \frac{1}{\sigma_n} \nabla_{\underline{w}} \sigma_n \quad (\text{chain rule})$$

$$= - \sum_{n=1}^N \frac{1}{\sigma_n} \cancel{\sigma_n} \cancel{(1-\sigma_n)} \nabla_{\underline{w}} \left. \begin{matrix} z^{(n)} \\ \underline{w}^T \underline{x}^{(n)} \end{matrix} \right\} \frac{\partial \sigma(a)}{\partial a} = \sigma(a)(1-\sigma(a))$$

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

