

# MLPR so far

Train  
val → splits, generalization, Gaussian stats, CLT  
Test

Pre-processing: look at data, one hot encoding, taking logs

"Generative model"

Bayes classifiers + Gaussians

Linear regression

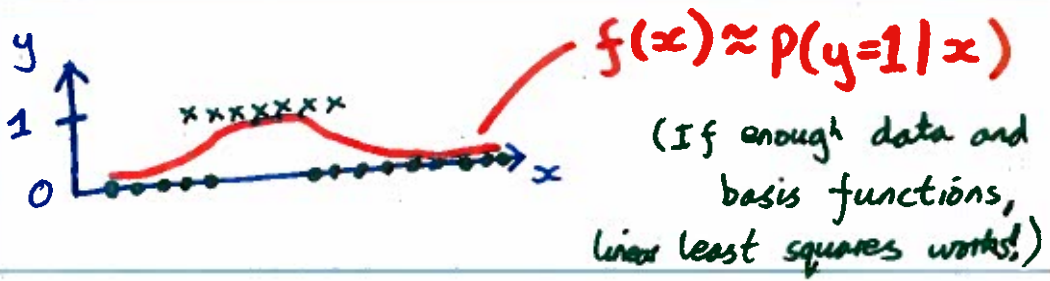
- with basis  $f_n$ 's
- regularization
- Bayesian prediction
- $\infty$  basis  $f_n$ 's: GPs

Today: revisit  
"Discriminative models"  
fit  $P(y|x)$  directly

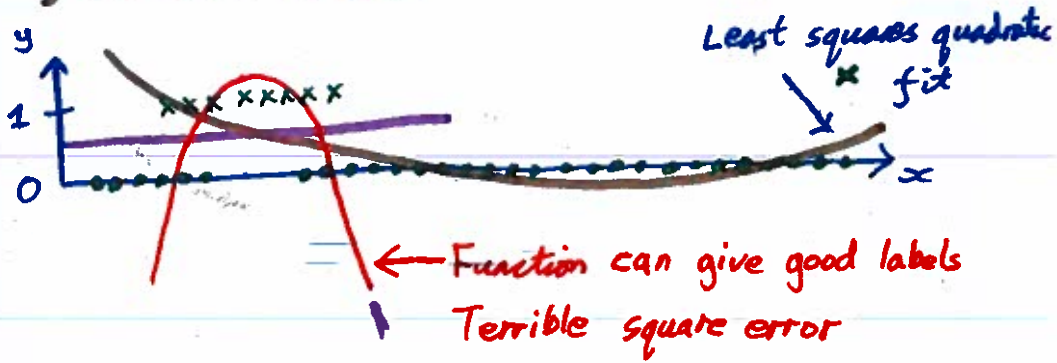
Start of non-linear  
models, fitted  
with gradient methods

Later will do Bayesian  
prediction with  
non-linear models

# Regressing on Labels



Often bad idea:



## Gradients for least squares cost

$$\text{Residuals } \underline{r} = \underset{N \times 1}{y} - \underset{f}{X \underline{w}}$$

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

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

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

"Gradient" vector of partial derivatives

$$\begin{aligned} \nabla_{\underline{w}} [r^T r] &= \underline{0} - 2 X^T y + 2 X^T X \underline{w} \\ &= -2 X^T \underbrace{(y - X \underline{w})}_{\underline{r}} \end{aligned}$$

### Gradient Descent

Initialize  $\underline{w}$  (to  $\underline{0}$ ?)

for  $t = 1 \dots T$ :

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

$\eta$  "eta", small number 0.01?

(Scratch working)

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

↑  
some  
vector

$$\frac{\partial \underline{w}^T \underline{h}}{\partial w_i} = \frac{\partial}{\partial w_i} \sum_j w_j h_j$$

$$= \frac{\partial}{\partial w_i} (w_1 h_1 + w_2 h_2 + \dots$$

$$\dots w_i h_i + \dots \\ \dots w_n h_n)$$

$$= h_i$$

Matrix Cookbook

# Normal Equations approach

$$\nabla_w [L^T L] = \underline{0} \quad \text{at least squares solution}$$

↳ all weights are happy where they are

$$(X^T X)w = X^T y$$

If  $(X^T X)^{-1}$  exists

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

Pseudo Inverse.

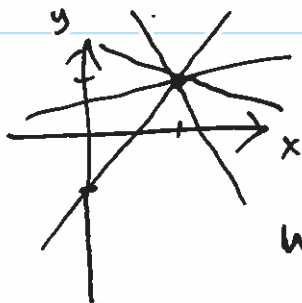
?  $X^{-1} (X^T X)^{-1} y$   $\rightarrow$  I

$= X^{-1} y$   $\times$

$N \times D$

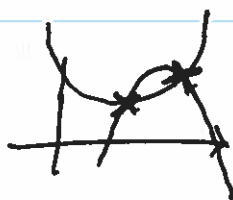
" $X \setminus y$ "

When is there a unique solution?



$$D=2$$

$$N=1$$



Different functions  
same train cost

We need  $N \geq D$

for unique solution.

One-hot encoding in R

"red"  $\rightarrow$  ~~1 0 0~~

"blue"  $\rightarrow$  ~~0 1 0~~

"green"  $\rightarrow$  ~~0 0 1~~  
 $x_1, x_2, x_3$

bias

1

1

1

$x_0$

Imagine I have weights  $\underline{w}$

New weights  $\underline{\tilde{w}}$ :

$$\tilde{w}_1 = w_1 + \delta$$

$$\tilde{w}_2 = w_2 + \delta$$

$$\tilde{w}_3 = w_3 + \delta$$

$$\tilde{w}_D = w_D - \delta$$

Same function  
 $\Rightarrow \underline{w}^T \underline{x} = \underline{\tilde{w}}^T \underline{x}$

for all  $\underline{x}$

(so same train cost)

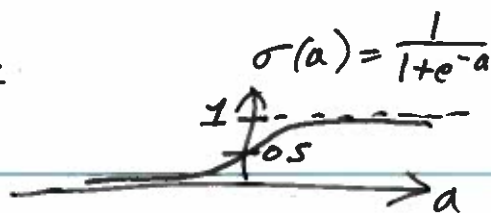
$$x_1 + x_2 + x_3 + x_D = 2$$

$$x_3 = 2 - x_1 - x_2 - x_D$$

# Logistic Regression

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

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



## Loss Function

Could use square loss again:

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

Interpretation:

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

Maximum Likelihood, maximize prob. of the data given  $\underline{w}$ :

$$P(\underline{y} | \underline{X}, \underline{w}) = \prod_{n=1}^N P(y^{(n)} | \underline{x}^{(n)}, \underline{w})$$

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:

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

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

$$\begin{aligned} \nabla_{\underline{w}} NLL &= - \sum_{n=1}^N \nabla_{\underline{w}} \log \sigma_n \\ &= - \sum_{n=1}^N \frac{1}{\sigma_n} \nabla_{\underline{w}} \sigma_n \quad (\text{chain rule}) \end{aligned}$$

$$\begin{aligned} &\frac{\partial \sigma(a)}{\partial a} = \sigma(a)(1 - \sigma(a)) \\ &= - \sum_{n=1}^N \frac{1}{\sigma_n} \sigma_n (1 - \sigma_n) \underbrace{\nabla_{\underline{w}} z^{(n)} \underline{w}^T \underline{x}^{(n)}}_{z^{(n)} \underline{x}^{(n)}} \\ &= - \sum_{n=1}^N (1 - \sigma_n) z^{(n)} \underline{x}^{(n)} \end{aligned}$$