# Bayes Classifiers

## Training time

Joint model $p(y, \underline{x}) = p(y) p(\underline{x}|y)$

$p(y=k) = \pi_k \approx \dfrac{\#\ k\ labels}{N}$

$p(\underline{x}|y=k) \ldots$ eg $N(\underline{x}; \mu^{(k)}, \Sigma^{(k)})$

Not Bayesian
Assuming we know all parameters.

Mean & cov of $\underline{x}$'s in class $k$

Naive Bayes $p(\underline{x}|y=k) = \prod_d p(x_d | y=k)$

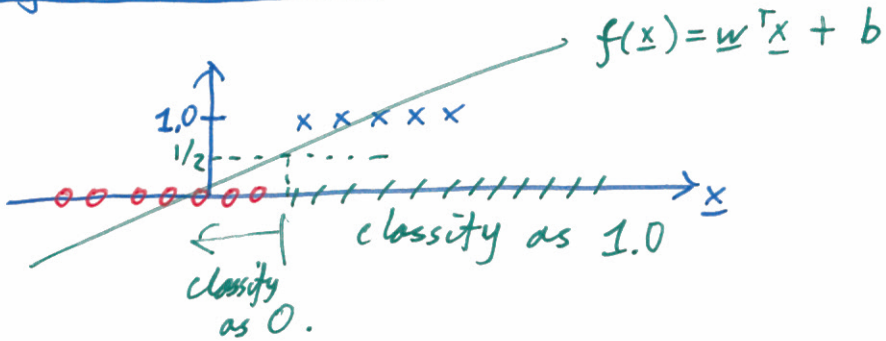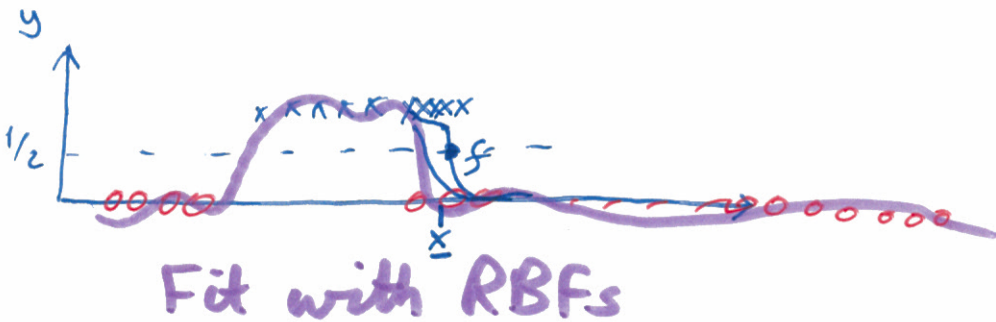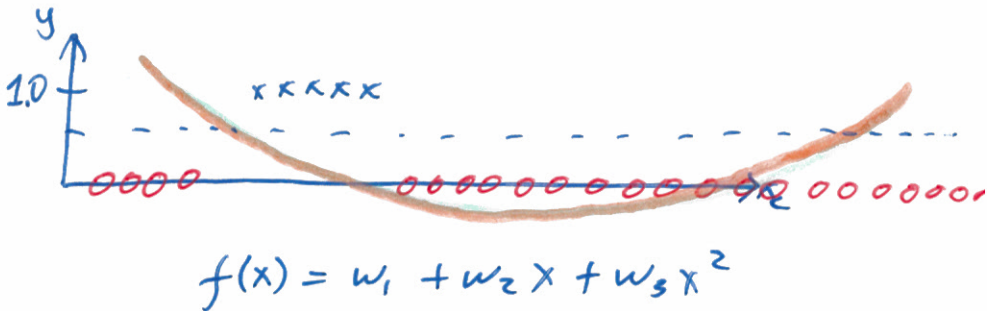Univariate Gaussian, Discrete, ...

## Test time

$p(y|\underline{x}) \propto p(y, \underline{x})$   (Bayes' Rule)

$y_{guess} = \underset{k}{argmax}\ p(y=k, \underline{x})$

"goodness"

# Regression to labels



$$f(\underline{x}) = \underline{w}^\top \underline{x} + b$$

classify as 1.0

classify as 0.

If $f(x) > \frac{1}{2}$, guess $y=1$



$$f(x) = w_1 + w_2 x + w_3 x^2$$



# Fit with RBFs

# If minimize square loss?

Minimize $\underset{P(y|\underline{x})}{E}\left[(y - f(\underline{x}))^2\right]$    at some location $\underline{x}$

Cost

$$= \underbrace{P_1}_{P(y=1|\underline{x})}(1-f)^2 + \underbrace{(1-P_1)}_{P(y=0|\underline{x})}(0-f)^2$$

$$= P_1(1-2f+f^2) + (1-P_1)f^2$$

$$= f^2(\cancel{P_1 - P_1}\ 1) - 2P_1 f + P_1$$
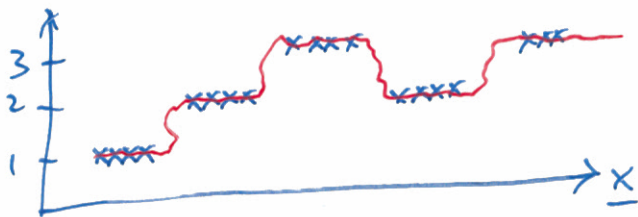
$$\frac{\partial cost}{\partial f} = 2f - 2P_1 = 0 \quad \underline{at\ optimum}$$

$$\boxed{f = P_1}$$

# Multiple classes.

$$y \in \{1, 2, 3, 4 \cdots 10\}$$

"sport"    "crime"    "romance"



If $\quad f(\underline{x}) = \underline{w}^T \underline{x}$

Maybe replace
$\underline{x}$ with $\phi(\underline{x})$

$f(\underline{x}^{(1)}) \approx 1 \Rightarrow$ "sport"

$f(\underline{x}^{(2)}) \approx 3 \Rightarrow$ "romance"

$f\left(\dfrac{\underline{x}^{(1)} + \underline{x}^{(2)}}{2}\right) \approx 2 \Rightarrow$ "crime"

# One-hot encoding, One-of-K encoding

Vector output

Is in $k^{th}$ position

$$y^{(n)} = [0 \ 00 \cdots 0 \ 1 \ 0 \cdots 0]^T$$

$K \times 1$ vector
If we have
K classes

If $n^{th}$ example
is in class $k$

## Fit K functions, one for each bit $y_k$

This pre-processing step is also useful
for input features

$$x_d \in \{\text{"red"}, \text{"green"}, \text{"blue"}\}$$
$$\{1, 2, 3\}$$

3 features    red $\rightarrow$   1   0   0    R doesn't
green $\rightarrow$   0   1   0    create
blue $\rightarrow$   0   0   1    this column.

Puzzle: In R you can do one-hot encoding

red $\rightarrow$ 10
green $\rightarrow$ 01
blue $\rightarrow$ 00

# Gradients for least squares cost

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

$\qquad\qquad\qquad \uparrow$ N×1 vector of scalar labels

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

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

"Gradient" vector of partial derivatives:

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

# Scratch working

$$\nabla_{\underline{w}} \left[ \underline{w}^T \underline{h} \right] = \begin{bmatrix} \dfrac{\partial \underline{w}^T \underline{h}}{\partial w_1} \\[2mm] \dfrac{\partial \underline{w}^T \underline{h}}{\partial w_2} \\[2mm] \cdot \\ \vdots \\ \dfrac{\partial \underline{w}^T \underline{h}}{\partial w_D} \end{bmatrix} = \underline{h}$$

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

$$= \frac{\partial}{\partial w_i} \left( w_1 h_1 + w_2 h_2 + \cdots \right.$$
$$\left. \underline{w_i h_i} + \cdots w_D h_D \right)$$

$$= h_i$$

Matrix Cookbook.