

BAYES CLASSIFIERS

At training time create model:

Labels: $P(y=k) = \pi_k$

Features: $P(\underline{x} | y=k)$ e.g. $\mathcal{N}(\underline{x}; \mu^{(k)}, \Sigma^{(k)})$

or discrete dist.

→ Models \underline{x} and y :

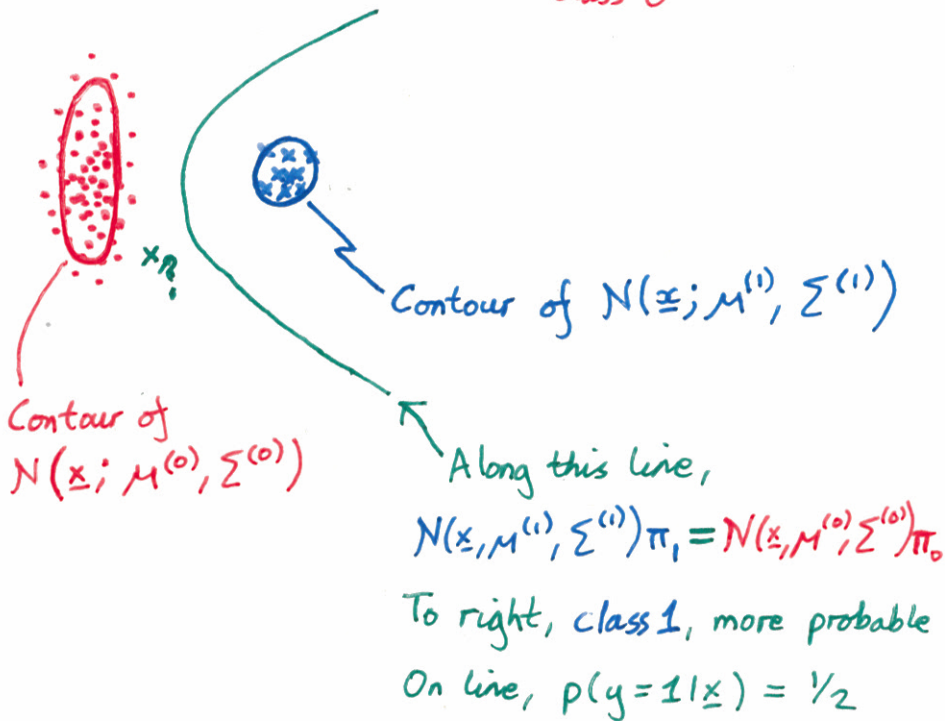
$$P(y, \underline{x}) = p(y) p(\underline{x} | y)$$

At test time use Bayes' rule:

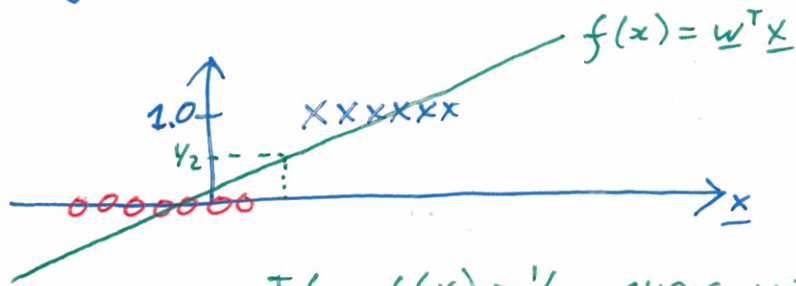
$$p(y | \underline{x}) \propto p(y, \underline{x})$$

$$p(y=k | \underline{x}) = \frac{p(y=k, \underline{x})}{\sum_{k'} p(y=k', \underline{x})}$$

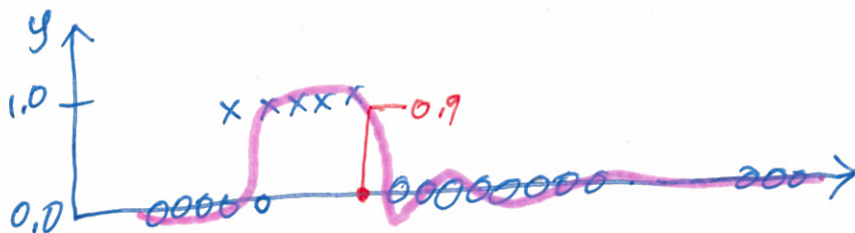
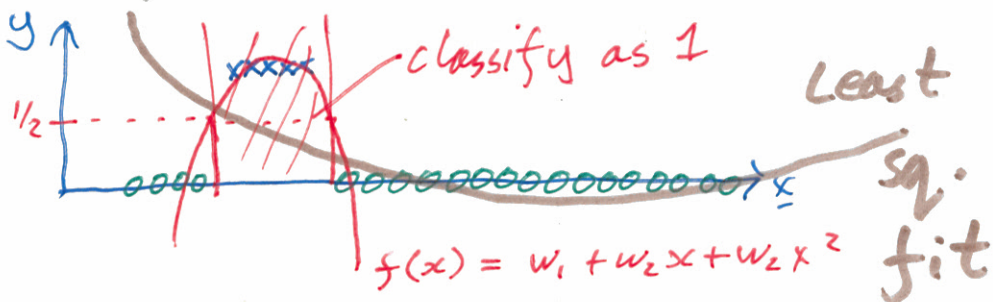
Labelled train points
• and * class 1
↖ class 0



Regressing on the labels



If $f(x) > 1/2$ guess $y = 1$



Radial basis f^n fit.

If minimize square loss

Minimize $\mathbb{E}_{p(y|x)}[(y-f)^2]$ at some location \underline{x}

$$= p_1(1-f)^2 + (1-p_1)(0-f)^2$$

↑

$p(y=1|x)$

↑

Two classes

$$p_0 + p_1 = 1$$

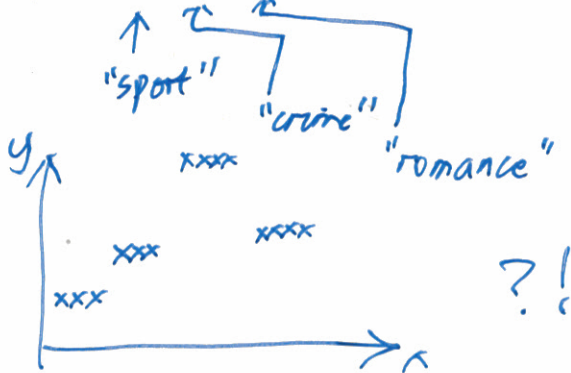
$$\text{cost} = f^2 - 2p_1 f + p_1$$

$$\frac{\partial \text{cost}}{\partial f} = 2f - 2p_1 = 0 \quad \text{at optimum}$$

$$f = p_1$$

Multiple classes

$$y = \{1, 2, 3, 4, \dots, 10\}$$



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

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

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

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

One-hot encoding, One-of-k encoding

Create vector output

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

↑
kth position

If example n is in class k .

Puzzle:

R comes with a f^n to do one-hot encoding
If you have k classes for linear regression

it returns a $k-1$ dim. vector

If you have class k it return $k-1$ zeros.

Gradients for least squares

$N \times D$ design matrix
of all inputs

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

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

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

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

"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:

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

η
step-size, small.