

rather than a hard rule, e.g. $\text{sign}(f(x, w))$
 what if we model $p(y_n = c_k | x_n)$

Approach 1: Discriminative modeling. We need to choose
 some model for $p(y_n = c_k | x_n)$

Binary case: sigmoid commonly used.

$$\sigma(z) = \frac{1}{1 + \exp\{-z\}}$$



put in a linear function $z = w^T x + w_0$

logistic regression (model for continuous probs)

turn the crank:

$$\hat{p}(y=1|z) = \frac{1}{1 + \exp\{-w^T x - w_0\}}$$

$$\begin{aligned} \hat{p}(y=0|x) &= 1 - \hat{p}(y=1|x) = \frac{\cancel{1} \exp\{-w^T x - w_0\}}{1 + \exp\{-w^T x - w_0\}} \\ &= \frac{1}{1 + \exp\{w^T x + w_0\}} \end{aligned}$$

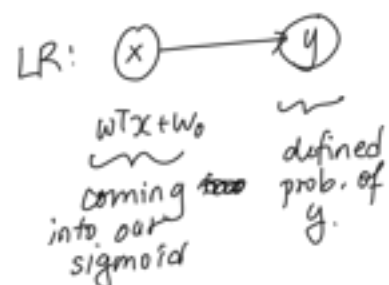
$$\text{lik } \ell(w) = \sum_{n=1}^N y_n \log \hat{p}(y=1|x_n) + \underbrace{(1-y_n) \log \hat{p}(y=0|x_n)}_{\text{goal: maximize}} \quad \hat{p}(1)^{y_n} \hat{p}(0)^{1-y_n}$$

$$\nabla \ell(w) = \sum_{n=1}^N -y_n x_n \hat{p}(y_n=0|x_n) + (1-y_n) x_n \hat{p}(y_n=1|x_n)$$

apply gradient descent/ascent

Note: very similar to the perceptron alg:

$$\begin{array}{lcl} \tilde{V}_P = \begin{cases} y_n=1 & x_n \\ y_n=0 & -x_n \end{cases} & \text{vs.} & \begin{array}{l} x_n \hat{p}(y_n=0) \\ -x_n \hat{p}(y_n=1) \end{array} \\ \text{if wrong} & & \text{recipe for SGD} \\ & \text{recipe for} & \text{w/ logit loss} \\ & \text{SGD w/} & \\ & \text{hinge loss} & \end{array}$$



\Rightarrow if I wanted to create a
 toy data set to test logistic
 regression.

\Rightarrow create a bunch of x 's

\Rightarrow choose w, w_0

\Rightarrow compute $\hat{p}(y|x)$

\Rightarrow sample $y \sim \hat{p}(y|x)$

Approach 2: Generative Model (model joint $p(x, y)$ rather than)

Approach 2: generative just $p(y|x)$

prior on $y \rightarrow (y) \rightarrow (x)$

- $\Rightarrow y \sim p(y)$
- \Rightarrow choose some parameters for $p(x|y)$
- \Rightarrow sample $x \sim p(x|y)$

if y is binary $p(y)$ can be a Bernoulli (easy)

if we want to infer $y|x$, we need Bayes rule:

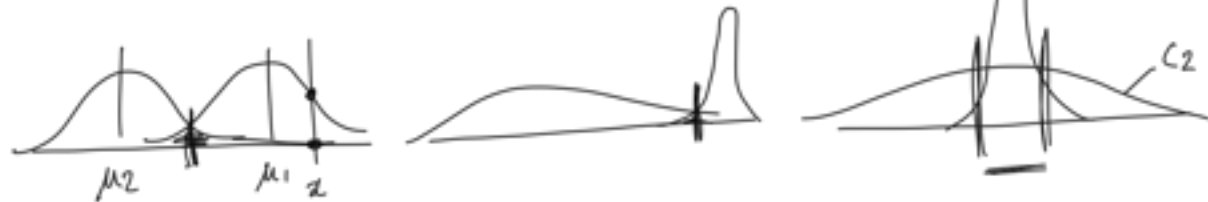
before regressed on this \swarrow

$$p(y=c_k|x) \propto \underbrace{p(x|y=c_k)}_{\text{form will depend on type of var. } x} \underbrace{p(y=c_k)}_{\text{prior (e.g. Bernoulli if binary)}} = p(x, y)$$

Suppose that x is continuous: can assume

$x \sim N(\mu_1, \Sigma_1)$ if $y \in C_1$ ✓

$x \sim N(\mu_2, \Sigma_2)$ if $y \in C_2$ ✓



(use $\log \frac{p(C_1|x)}{p(C_2|x)}$ as the boundary

$$\frac{1}{\sqrt{2\pi}|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right\}$$

$$\log \left(\frac{p(C_1)(2\pi)^{d/2} |\Sigma_2|^{1/2}}{p(C_2)(2\pi)^{d/2} |\Sigma_1|^{1/2}} \right) - \exp\left\{-\frac{1}{2} \left[(x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1) - (x-\mu_2)^T \Sigma_2^{-1} (x-\mu_2) \right]\right\}$$

= a lot of terms that don't depend on x $- \frac{1}{2} \left[(x)(\Sigma_1^{-1} - \Sigma_2^{-1})(x)^T - 2x(\Sigma_1^{-1}\mu_1^T + \Sigma_2^{-1}\mu_2^T) + \text{more terms without } x \right]$

= a lot of terms w/o x + $x(\Sigma_1^{-1} - \Sigma_2^{-1})x^T + x(\text{something})$

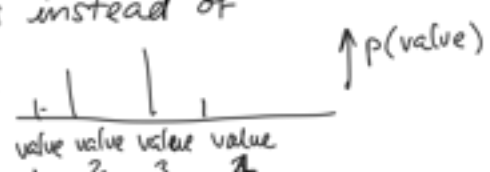
if $\Sigma_1 = \Sigma_2$
then we have
a linear boundary!



That was the x is continuous case...

$p(y \in C_1|x)$ $p(y \in C_2|x)$

discrete case: now we can multinomials instead of Gaussians to define $p(x|y)$. Multinomial:



→ explode w/ number of dims (R^D)

make an assumption! $p(x|c) = \prod_{d=1}^D p(x_d|c)$



Naive Bayes assumption

We can write down likelihood: $\pi_{dj}^k = \Pr(x_d \text{ takes value } j \text{ in } y=c_k)$

$$p(x|\pi, y) = \prod_{n=1}^N \prod_{d=1}^D \pi_{dj}^k \quad \text{where } y_n = k \text{ and } x_{nd} = j$$

$$\log p(x|\pi, y) = \sum_n \sum_d \log \pi_{dj}^k \mathbb{I}(y_n=k) \mathbb{I}(x_{nd}=j) = \sum_n \sum_d \sum_k \log \pi_{dj}^k \mathbb{I}(y_n=k) \mathbb{I}(x_{nd}=j)$$

tip: group by class.

ASIDE: / NOTES:

1) we can do full Bayes (e.g. priors over params) if we wish

2) we can do multi-class easily for generative case

(a little harder in the discrim case: separate w_k for each class...)



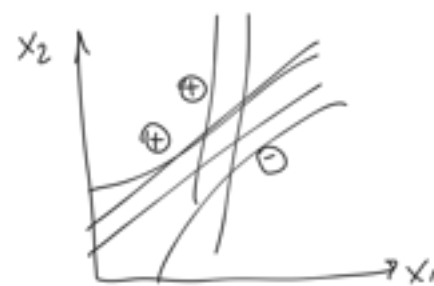
3) generative model: easy to deal w/ missing data

$$\underbrace{p(x|y) p(y)}_{\text{Bayes Rule}} = \left[\prod_{d=\text{we observe}} p(x_d|y) \right] \cdot p(y)$$

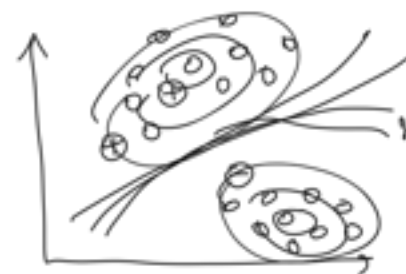
$$\frac{\sigma(w^T x + w_0)}{p(y|x)}$$

Concept Check: Semi-Supervised Learning

How can unlabeled data help?



What boundaries might you fit?



does unlabeled help?

... are they both as easy to use

generative vs. discriminative models: are they both
in this semi-supervised setting? why or why not?

assume: density of x reflects label y



y
 \downarrow
 x

$$\frac{p(y|x)}{1}$$