

## Bayes Decision Rule

Example:  $x$  describes patient,  $y$  whether drug will cure patient

know  $P(y=1|x=x) = 0.7$

$P(y=-1|x=x) = 0.3$

predicting  $y=1$ : make error with probability 0.3

$y=-1$ : make error with probability 0.7

→ predict label  $y$  that has the highest  $P(y|x=x)$

## Bayes Error Rate

What is  $err_P(h_{\text{Bayes}})$  for 0/1 loss?

Given an instance  $x$ ,  $h_{\text{Bayes}}$  has probability of error of  $\min_{y \in \mathcal{Y}} (1 - P(y=y|x=x))$

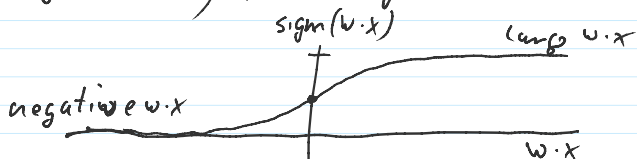
→ Expectation over all  $x \sim P(x)$ :

$$err_P(h_{\text{Bayes}}) = \sum_{x \in \mathcal{X}} P(x=x) \min_{y \in \mathcal{Y}} (1 - P(y=y|x=x))$$

## Logistic Regression

Idea: Model  $P(y|x) = \text{Ber}(y|\mu = \text{sigm}(w \cdot x))$  for binary classification

Sigmoid:  $\text{sigm}(w \cdot x) = \frac{1}{1 + e^{-w \cdot x}}$



$P(y=1|x, w) = \text{sigm}(w \cdot x) = \frac{1}{1 + e^{-w \cdot x}}$

$P(y=-1|x, w) = 1 - \text{sigm}(w \cdot x) = 1 - \frac{1}{1 + e^{-w \cdot x}} = \frac{1 + e^{-w \cdot x} - 1}{1 + e^{-w \cdot x}} = \frac{e^{-w \cdot x}}{1 + e^{-w \cdot x}}$

→  $P(y|x, w) = \frac{1}{1 + e^{-y w \cdot x}} = \frac{1}{1 + e^{w \cdot x}}$

Training: Conditional Maximum Likelihood

$$w = \underset{w}{\text{argmax}} P(y_1 \dots y_m | x_1 \dots x_m, w)$$

$$= \underset{w}{\text{argmax}} \prod_{i=1}^m P(y_i | x_i, w)$$

$$\begin{aligned}
&= \arg \min_w \prod_{i=1}^m (1 + e^{-y_i w \cdot x_i}) \\
&= \arg \max_w \prod_{i=1}^m \frac{1}{1 + e^{-y_i w \cdot x_i}} \\
&= \arg \min_w - \log \left[ \prod_{i=1}^m \frac{1}{1 + e^{-y_i w \cdot x_i}} \right] \\
&= \arg \min_w \sum_{i=1}^m - \log \left( \frac{1}{1 + e^{-y_i w \cdot x_i}} \right) \\
&= \arg \min_w \sum_{i=1}^m \log (1 + e^{-y_i w \cdot x_i})
\end{aligned}$$

### Regularized Logistic Regression

Idea: Like Logistic Regression, but put prior  $P(w)$  on the parameters

$$\rightarrow P(w) = N(w | 0, \text{diag}(\sigma)) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{w \cdot w}{2\sigma^2}}$$

Training: Maximum a posteriori estimate

$$\begin{aligned}
w &= \arg \max_w P(y_1 \dots y_m | x_1 \dots x_m, w) \cdot P(w) \\
&= \arg \min_w \left[ \sum_{i=1}^m \log (1 + e^{-y_i w \cdot x_i}) \right] - \log \left( \frac{1}{\sqrt{2\pi}\sigma} \right) - \log \left( e^{-\frac{w \cdot w}{2\sigma^2}} \right) \\
&= \arg \min_w \frac{1}{2} w \cdot w + \sigma^2 \sum_{i=1}^m \log (1 + e^{-y_i w \cdot x_i})
\end{aligned}$$

$$\rightarrow \text{SVM Primal: } w = \arg \min_w \frac{1}{2} w \cdot w + C \sum_{i=1}^m \max(0, 1 - y_i w \cdot x_i)$$

### Ridge Regression

Regression: Predict real-valued  $y \in \mathbb{R}$  from features  $x \in \mathbb{R}^U$

$$\text{Idea: Model likelihood } P(y | x, w) = N(y | w \cdot x, \epsilon) = \frac{1}{\sqrt{2\pi}\epsilon} e^{-\frac{(w \cdot x - y)^2}{2\epsilon^2}}$$

$$\text{and prior } P(w) = N(w | 0, \text{diag}(\sigma))$$

Training: Max a posteriori

$$\begin{aligned}
w &= \arg \max_w \left[ \prod_{i=1}^m P(y_i | x_i, w) \right] \cdot P(w) \\
&= \arg \min_w - \left[ \sum_{i=1}^m \log P(y_i | x_i, w) \right] - \log P(w) \\
&= \arg \min_w \left[ \frac{1}{2} \sum_{i=1}^m (w \cdot x_i - y_i)^2 + \dots \right]
\end{aligned}$$

$$\begin{aligned}
 &= \arg \min_w \left[ \frac{1}{2\epsilon^2} \sum_{i=1}^m (w \cdot x_i - y_i)^2 \right] - \frac{1}{2\alpha^2} w \cdot w \\
 &= \arg \min_w \frac{1}{2} w \cdot w + \frac{\sigma^2}{2\epsilon^2} \sum_{i=1}^m (w \cdot x_i - y_i)^2
 \end{aligned}$$

## Regularized (linear) Models

L:

- Hinge loss:  $\max(0, 1 - y\hat{y})$

Exp loss:  $\exp(-\beta y \cdot \hat{y})$

Squared error:  $(y - \hat{y})^2$

Absolute error:  $|y - \hat{y}|$

R:

$L_2$  norm:  $\|w\|_2 = \sqrt{w \cdot w}$

$L_1$  norm:  $\|w\|_1 = \sum |w_i|$

$L_0$  norm:  $\|w\|_0$

$L_p$  norm

Beyond linear models:

Linear:  $w \cdot x + b$

Deep network:  $f_w(x)$

Kernel:  $w \cdot \phi(x) + b$