

SVMs and Kernels



Two different approaches to regression/classification

Generative Approach

$$x \sim p \quad y = w^T x + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$p(y|x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(y - w^T x)^2}{2\sigma^2}\right)$$

- Assume something about $P(x, y)$
- Find f which maximizes likelihood of training data | assumption $\{(x_i, y_i)\}_{i=1}^n$ w maximize $\prod_{i=1}^n p(x_i, y_i | w)$

- Often reformulated as minimizing loss

$$\Leftrightarrow \min_w -\sum_{i=1}^n \log p(x_i, y_i | w)$$

Versus

Discriminative Approach

- Pick a loss function
- Pick a set of hypotheses H
- Pick f from H which minimizes loss on training data

l_2 : regression

0/1: classification

linear function

$$f(x) = w^T x$$

$$\begin{cases} w^T x + b > 0 \Rightarrow 1 \\ w^T x + b < 0 \Rightarrow 0 \end{cases}$$

$$\min_{f \in H} \frac{1}{n} \sum_{i=1}^n l(f(x_i), y_i)$$

★ No assumption on $p(y|x)$

Our description of logistic regression was the former

- Learn: $f: \mathbf{X} \rightarrow \mathbf{Y}$

- \mathbf{X} – features

- \mathbf{Y} – target classes

- Loss function:

$$l(f(x), y) = \mathbb{1}\{f(x) \neq y\}$$

- Expected loss of f : $Y \in \{-1, 1\}$

$$\mathbb{E}_{X,Y} [\mathbb{1}\{f(X) \neq Y\}] = \mathbb{E}_X [\mathbb{E}_{Y|X} [\mathbb{1}\{f(X) \neq Y\} | X=x]]$$

$$\mathbb{E}_{Y|X} [\mathbb{1}\{f(X) \neq Y\} | X=x] = 1 - P(Y = f(X) | X=x)$$

- Bayes optimal classifier: $f(x) = \underset{y}{\operatorname{argmax}} \underbrace{P(Y=y | X=x)}_{\text{need a model}}$

- Model of logistic regression:

$$P(Y=y | X, w) = \frac{1}{1 + \exp(-y w^T x)}$$

Our description of logistic regression was the former

- Learn: $f: \mathbf{X} \rightarrow \mathbf{Y}$

- \mathbf{X} – features

- \mathbf{Y} – target classes

$$Y \in \{-1, 1\}$$

- Expected loss of f :

$$\mathbb{E}_{XY}[\mathbf{1}\{f(X) \neq Y\}] = \mathbb{E}_X[\mathbb{E}_{Y|X}[\mathbf{1}\{f(x) \neq Y\}|X = x]]$$

$$\mathbb{E}_{Y|X}[\mathbf{1}\{f(x) \neq Y\}|X = x] = 1 - P(Y = f(x)|X = x)$$

- Bayes optimal classifier:

$$f(x) = \arg \max_y \mathbb{P}(Y = y|X = x)$$

- Model of logistic regression:

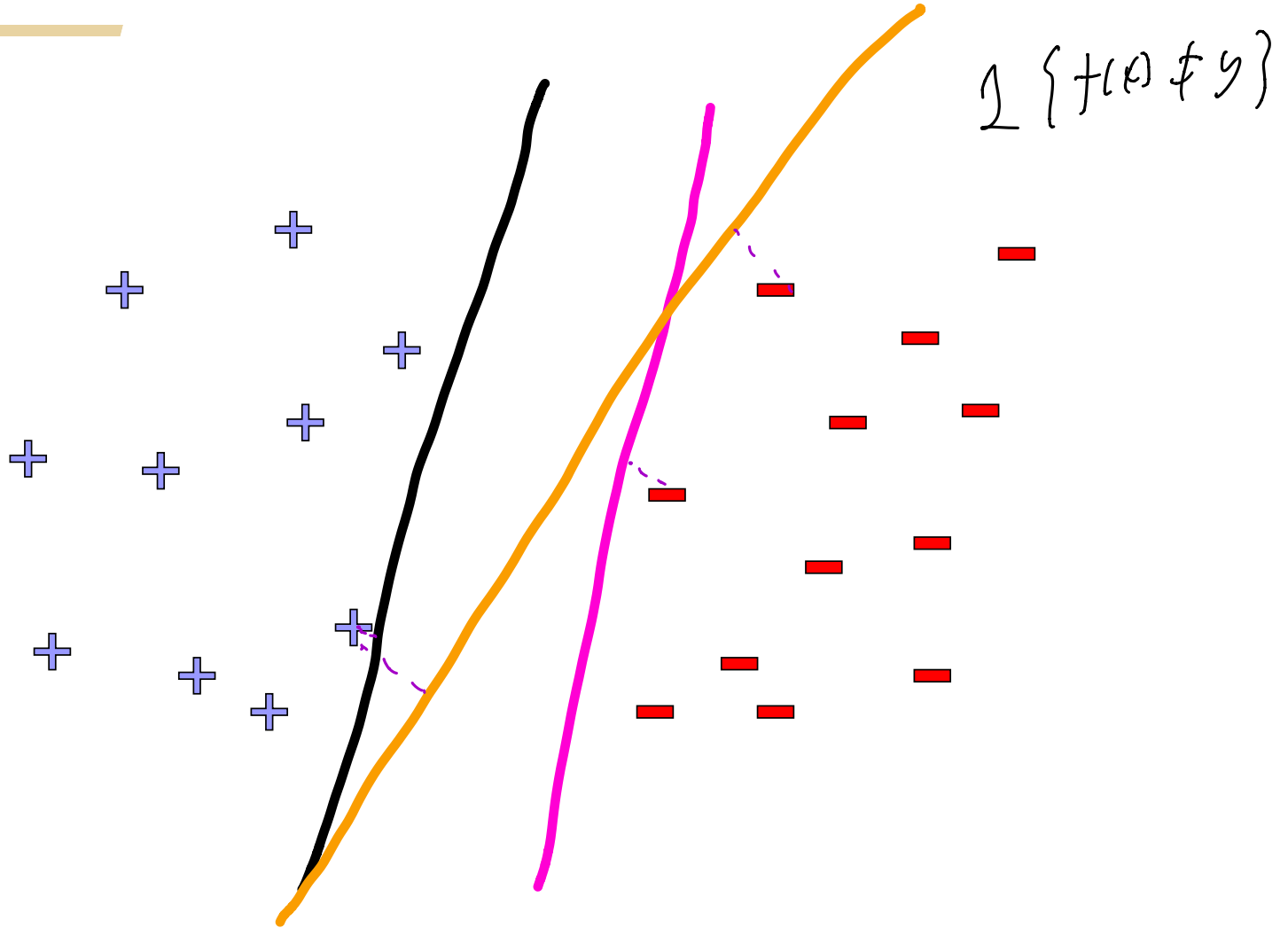
generative

$$P(Y = y|x, w) = \frac{1}{1 + \exp(-y w^T x)}$$

discriminative

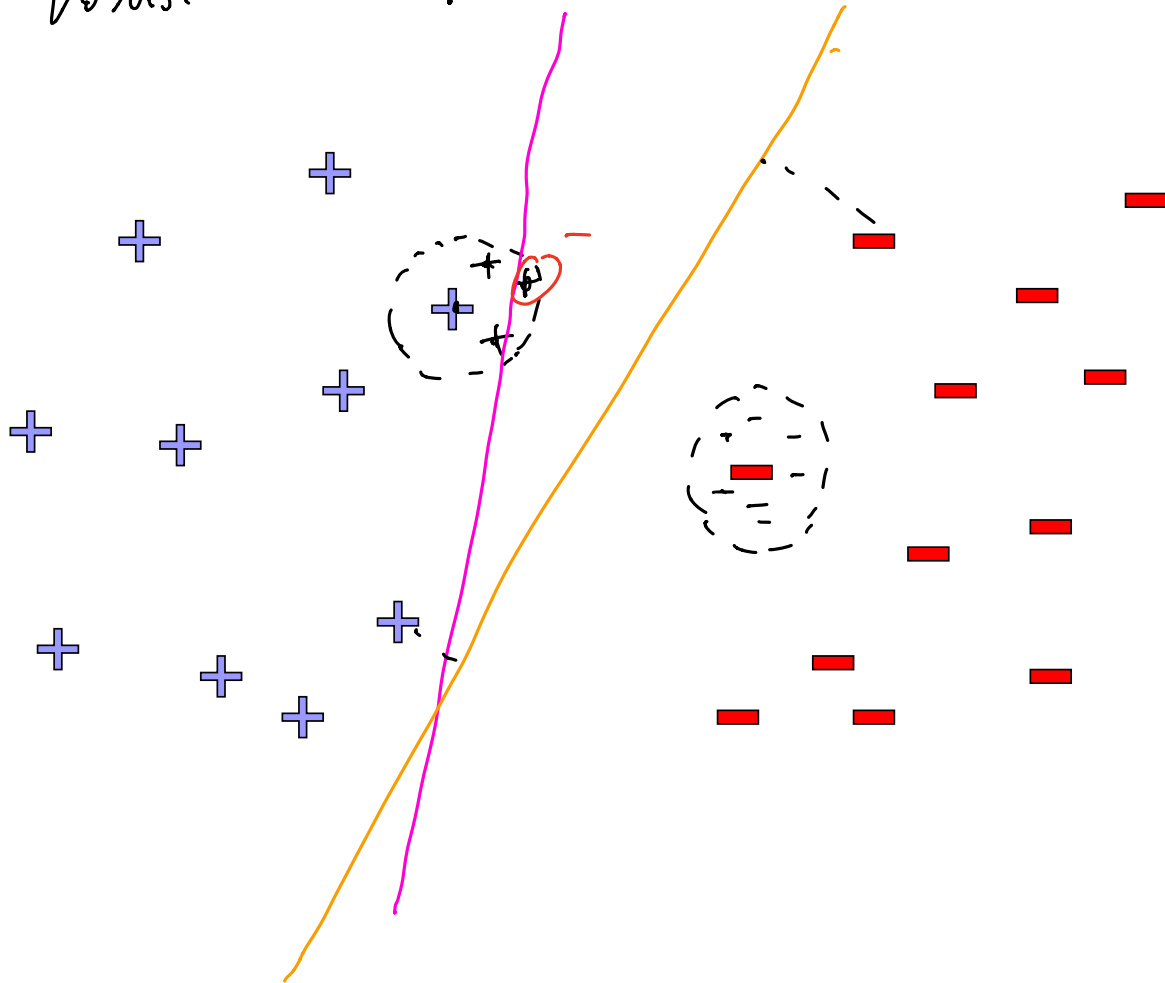
What if the model is wrong? What other ways can we pick linear decision rules?

Linear classifiers – Which line is better?



Linear classifiers – Which line is better?

robustness to perturbation



Linear classifiers – Which line is better?

