

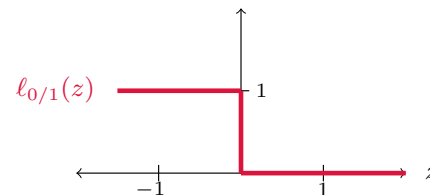
## Objectives

## Prediction error / zero-one loss

$P$  is a distribution over  $\mathcal{X} \times \{-1, +1\}$ , and  $(X, Y) \sim P$ .

For any classifier  $f: \mathcal{X} \rightarrow \{-1, +1\}$ ,

$$\text{err}(f) = P(f(X) \neq Y) = \mathbb{E}[\ell_{0/1}(Yf(X))].$$



Also works with **real-valued predictors**  $f: \mathcal{X} \rightarrow \mathbb{R}$ ; for example:

- ▶ **k-NN**: average of  $y$ -values of  $k$  nearest neighbors.
- ▶ **Trees**: leaf nodes with a real-valued output (e.g., average of  $y$ -values of training examples that reach a leaf). “Regression trees”
- ▶ **Linear classifiers**:  $x \mapsto \langle w, x \rangle - t$ .
- ▶ **Classifiers from generative models**:  $x \mapsto P_{\hat{\theta}}(Y = +1 | X = x) - 1/2$ .

Often useful to adjust threshold (e.g.,  $t$  and  $1/2$  above).

1 / 11

2 / 11

## Thresholds

Uses for adjusting threshold  $t$

Often have **different costs for different kinds of mistakes**:

	$f(X) \leq t$	$f(X) > t$
$Y = -1$	0	$c$
$Y = +1$	$1 - c$	0

Also, often interested in **different performance criteria**.

▶ **Precision**:

$$P(Y = +1 | f(X) > t)$$

▶ **Recall** (a.k.a. **Sensitivity**, **True Positive Rate**):

$$P(f(X) > t | Y = +1)$$

▶ **Specificity**:

$$P(f(X) \leq t | Y = -1)$$

▶ **False Positive Rate**:

$$P(f(X) > t | Y = -1)$$

3 / 11

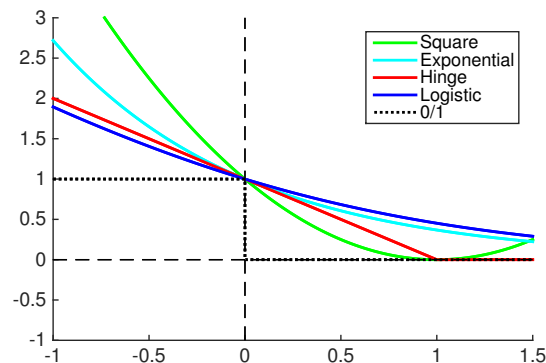
## Conditional probability estimation

Sometimes would like real-valued predictor  $f$  to be related to the **conditional probability function**  $\eta$

$$\eta(x) = P(Y = +1 | X = x).$$

- ▶ Straightforward when using generative models.
- ▶ Can use a loss function that is minimized by  $\eta$  (or some invertible transformation thereof).

4 / 11



**Goal:** loss function that is minimized by (some invertible transformation of) the conditional probability function

$$\eta(x) = P(Y = +1 \mid X = x).$$

Loss functions and their minimizers

- **Square loss:**  $\ell_{\text{sq}}(z) = (1 - z)^2$   
 $\mathbb{E}[\ell_{\text{sq}}(Y f(x)) \mid X = x]$  is minimized by  $f$  s.t.  $f(x) = 2\eta(x) - 1$ .  
So  $\eta(x) = (f(x) + 1)/2$ .
- **Logistic loss:**  $\ell_{\text{logistic}}(z) = \ln(1 + \exp(-z))$   
 $\mathbb{E}[\ell_{\text{logistic}}(Y f(x)) \mid X = x]$  is minimized by  $f$  s.t.  $f(x) = \ln\left(\frac{\eta(x)}{1 - \eta(x)}\right)$ .  
So  $\eta(x) = (1 + \exp(-f(x)))^{-1}$ .

Non-example

- **Hinge loss:**  $\ell_{\text{hinge}}(z) = \max\{0, 1 - z\}$   
 $\mathbb{E}[\ell_{\text{hinge}}(Y f(x)) \mid X = x]$  is minimized by  $f$  s.t.  $f(x) = \text{sign}(2\eta(x) - 1)$ .  
Cannot recover  $\eta(x)$  from  $f(x)$ .

**Using loss functions:** easy with linear/affine functions whenever the loss function  $\ell$  is a convex function:

$$\min_{\mathbf{w} \in \mathbb{R}^d} R(\mathbf{w}) + \frac{1}{n} \sum_{i=1}^n \ell(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle).$$

(Here, the regularization function  $R$  is also assumed to be convex.)

**Caveat:** Might not be possible to represent

$$x \mapsto 2\eta(x) - 1 \quad \text{or} \quad x \mapsto \ln\left(\frac{\eta(x)}{1 - \eta(x)}\right)$$

as (say) a linear function  $x \mapsto \langle \mathbf{w}, x \rangle$ .

**Common remedies:** enhance the feature space via feature expansion or kernels, or use more flexible models (e.g., tree models).

Sometimes  $\mathcal{Y}$  is not just  $\{0, 1\}$  or  $\{1, 2, \dots, K\}$ , but rather a collection of *structured objects*.

Example: sequence tagging

- $\mathcal{X}$ : sequences of English words
- $\mathcal{Y}$ : sequences of parts-of-speech

the/D man/N saw/V the/D dog/N

(Verbs tend to follow Nouns.)

Many other examples:

- sentence parse trees
- web search result ranking
- visual scene labeling
- ...

Featurization

Create several input-output feature maps  $\phi_1, \phi_2, \dots, \phi_d: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ .

- ▶ e.g.,  $\phi_{1000}(\mathbf{x}, \mathbf{y}) = \mathbb{1}\{i\text{-th word in } \mathbf{x} \text{ is "the", and } i\text{-th POS in } \mathbf{y} \text{ is "D"}\}$

For each possible  $\mathbf{y} \in \mathcal{Y}$ , consider an *input-output feature vector*:

$$\Phi(\mathbf{x}, \mathbf{y}) := (\phi_1(\mathbf{x}, \mathbf{y}), \phi_2(\mathbf{x}, \mathbf{y}), \dots, \phi_d(\mathbf{x}, \mathbf{y})) \in \mathbb{R}^d.$$

**Note:** often  $d$  is enormous, but  $\phi_i(\mathbf{x}, \mathbf{y}) = 0$  for most  $i$ .

Model

Prediction model is based on *linear functions of input-output feature vectors*:

$$\mathbf{x} \mapsto \arg \max_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}) \rangle$$

for weight vector  $\mathbf{w} \in \mathbb{R}^d$ .

**Note:** the  $\arg \max$  can often be computed efficiently (e.g., via dynamic programming), even when  $\mathcal{Y}$  is enormous.

Online Structured Perceptron

**input** Labeled examples  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$  from  $\mathcal{X} \times \mathcal{Y}$ .

- 1: **initialize**  $\hat{\mathbf{w}}_1 := \mathbf{0}$ .
- 2: **for**  $t = 1, 2, \dots$ , **do**
- 3:   Predict:  $\hat{\mathbf{y}}_t := \arg \max_{\mathbf{y} \in \mathcal{Y}} \langle \hat{\mathbf{w}}_{t-1}, \Phi(\mathbf{x}_t, \mathbf{y}) \rangle$
- 4:   **if**  $\hat{\mathbf{y}}_t \neq \mathbf{y}_t$  **then**
- 5:     Update:

$$\hat{\mathbf{w}}_t := \hat{\mathbf{w}}_{t-1} + \Phi(\mathbf{x}_t, \mathbf{y}_t) - \Phi(\mathbf{x}_t, \hat{\mathbf{y}}_t).$$

- 6:   **else**
- 7:     No update:  $\hat{\mathbf{w}}_t := \hat{\mathbf{w}}_{t-1}$
- 8:   **end if**
- 9: **end for**

Can also help to make multiple passes through data, and also to employ averaging (as in Averaged Perceptron).

Key takeaways

- 1. Concept of real-valued predictors and thresholds; alternative performance criteria.
- 2. Eliciting conditional probabilities with loss functions.
- 3. High-level idea of structured output prediction.