

Notation

$\mathcal{X}$  :  $p$ -dim. **input space**

Usually we assume  $\mathcal{X} = \mathbb{R}^p$ , but categorical **features** can also occur

$\mathcal{Y}$  : **target space**

E.g.:  $\mathcal{Y} = \mathbb{R}$ ,  $\mathcal{Y} = \{0, 1\}$ ,  $\mathcal{Y} = \{-1, 1\}$ ,  $\mathcal{Y} = \{1, \dots, g\}$

$\mathbf{x} = (x_1, \dots, x_p)^T \in \mathcal{X}$  : **feature vector**

$y \in \mathcal{Y}$  : **target / label / output**

$\mathbb{P}_{xy}$  : **Joint probability distribution on  $\mathcal{X} \times \mathcal{Y}$**

$p(\mathbf{x}, y)$  or  $p(\mathbf{x}, y \mid \theta)$  : **joint probability density function (pdf)**

$(\mathbf{x}^{(i)}, y^{(i)})$  :  $i$  -th **observation** or **instance**

$\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})\}$

**data set** with  $n$  observations.  
 $\mathcal{D}_{\text{train}}, \mathcal{D}_{\text{test}}$  : data for training and testing

Often,  $\mathcal{D} = \mathcal{D}_{\text{train}} \dot{\cup} \mathcal{D}_{\text{test}}$ .

$f(\mathbf{x})$  or  $f(\mathbf{x} \mid \theta) \in \mathbb{R}$  or  $\mathbb{R}^g$  : prediction function (**model**)

We might suppress  $\theta$  in notation.

$h(\mathbf{x})$  or  $h(\mathbf{x} \mid \theta) \in \mathcal{Y}$

Discrete prediction for classification.

$\theta = (\theta_1, \theta_2, \dots, \theta_d) \in \Theta$ : model **parameters**

Some models may traditionally use different symbols.

$\mathcal{H}$  : **hypothesis space**

$f$  lives here, restricts the functional form of  $f$ .

$\epsilon = y - f(\mathbf{x})$  or  $\epsilon^{(i)} = y^{(i)} - f(\mathbf{x}^{(i)})$

**Residual** in regression.

$yf(\mathbf{x})$  or  $y^{(i)}f(\mathbf{x}^{(i)})$  : **margin** for binary classification

With,  $\mathcal{Y} = \{-1, 1\}$ .

$\pi_k(\mathbf{x}) = \mathbb{P}(y = k \mid \mathbf{x})$ : **posterior probability** for class  $k$ , given  $\mathbf{x}$

In case of binary labels we might abbreviate  $\pi(\mathbf{x}) = \mathbb{P}(y = 1 \mid \mathbf{x})$ .

$\pi_k = \mathbb{P}(y = k)$ : **prior probability** for class  $k$

In case of binary labels we might abbreviate  $\pi = \mathbb{P}(y = 1)$ .

$\mathcal{L}(\theta)$  and  $\ell(\theta)$  : Likelihood and log-Likelihood for a parameter  $\theta$

These are based on a statistical model.

$\hat{y}, \hat{f}, \hat{h}, \hat{\pi}_k(\mathbf{x}), \hat{\pi}(\mathbf{x})$  and  $\hat{\theta}$

These are learned functions and parameters ( These are estimators of corresponding functions and parameters).

Concepts

**Model:**  $f : \mathcal{X} \rightarrow \mathbb{R}^g$  is a function that maps feature vectors to predictions.

**Learner:** takes a data set with features and outputs (**training set**,  $\in \mathcal{X} \times \mathcal{Y}$ ) and produces a **model** (which is a function  $f : \mathcal{X} \rightarrow \mathbb{R}^g$ )

**Learning = Representation + Evaluation + Optimization.**

**Representation:** (Hypothesis space) Defines which kind of model structure of  $f$  can be learned from the data.  
Example: Linear functions, Decision trees etc.

**Evaluation:** A metric that quantifies how well a specific model performs on a given data set. Allows us to rank candidate models in order to choose the best one.  
Example: Squared error, Likelihood etc.

**Optimization:** Efficiently searches the hypothesis space for good models.  
Example: Gradient descent, Quadratic programming etc.

**Loss function:** The “goodness” of a prediction  $f(\mathbf{x})$  is measured by a loss function  $L(y, f(\mathbf{x}))$

Through **loss**, we calculate the prediction error and the choice of the loss has a major influence on the final model

**Risk Minimization:** The ability of a model  $f$  to reproduce the association between  $\mathbf{x}$  and  $y$  that is present in the data  $\mathcal{D}$  can be measured by the average loss: the **empirical risk**.

$$\mathcal{R}_{\text{emp}}(f) = \frac{1}{n} \sum_{i=1}^n L\left(y^{(i)}, f\left(\mathbf{x}^{(i)}\right)\right).$$

Learning then amounts to **empirical risk minimization** – figuring out which model  $f$  has the smallest average loss:

$$\hat{f} = \arg \min_{f \in \mathcal{H}} \mathcal{R}_{\text{emp}}(f).$$

Regression Losses

**Basic Idea (L2 loss/ squared error):**

- $L(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$  or  $L(y, f(\mathbf{x})) = 0.5(y - f(\mathbf{x}))^2$
- Convex and differentiable.
- Tries to reduce large residuals (if residual is twice as large, loss is 4 times as large)

**Basic Idea (L1 loss/ absolute error):**

- $L(y, f(\mathbf{x})) = |y - f(\mathbf{x})|$
- Convex and more robust
- No derivatives for  $= 0$ ,  $y = f(\mathbf{x})$ , optimization becomes harder
- $\hat{f}(\mathbf{x}) = \text{median of } y \mid \mathbf{x}$

Classification

We want to assign new observations to known categories according to criteria learned from a training set.

Assume we are given a **classification problem**:

$$\begin{array}{ll} \mathbf{x} \in \mathcal{X} & \text{feature vector} \\ y \in \mathcal{Y} = \{1, \dots, g\} & \text{categorical output variable (label)} \\ \mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})\} & \text{observations of } \mathbf{x} \text{ and } y \end{array}$$

Classification usually means to construct  $g$  discriminant functions:  $f_1(\mathbf{x}), \dots, f_g(\mathbf{x})$ , so that we choose our class as  $h(\mathbf{x}) = \arg \max_k f_k(\mathbf{x})$  for  $k = 1, 2, \dots, g$

**Linear Classifier:**

If the functions  $f_k(\mathbf{x})$  can be specified as linear functions, we will call the classifier a *linear classifier*.

**Note:** All linear classifiers can represent non-linear decision boundaries in our original input space if we include derived features. For example: higher order interactions, polynomials or other transformations of  $\mathbf{x}$  in the model.

**Binary classification:** If only 2 classes exist

We can use a single discriminant function  $f(\mathbf{x}) = f_1(\mathbf{x}) - f_2(\mathbf{x})$ .