# Lecture 1    Mathematical

# Foundations of Machine Learning

# Elements of Machine Learning

Imagine we want to take a photograph of a face and determine whether the person is smiling or not.

Key idea: we represent whether face is smiling with a <u>model</u> - a mathematical description of the data

Basis steps:

① Collect raw data — e.g. photographs of faces

② preprocessing — change the data to simplify subsequent operations without losing relevant information — e.g. crop images to only contain <u>one</u> face

Center face

resize so all images have same # pixels

③ Feature extraction — reduce raw data by extracting features or properties relevant to model (This step if often unnecessary in modern image recognition systems based on deep neural nets)
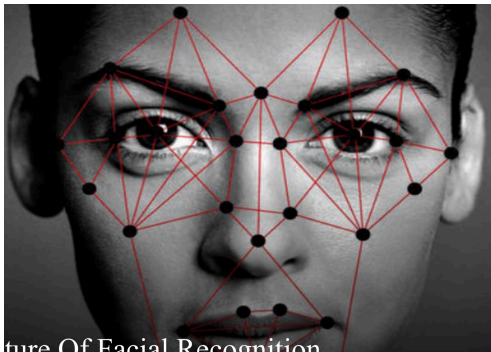

ture Of Facial Recognition

- e.g. distances between pairs of facial landmarks

④ generate *training samples* = large collection of examples we can use to learn a model

$$(\underline{x}_i, y_i) \text{ for } i=1,\ldots,n \Rightarrow n = \text{\# training samples}$$
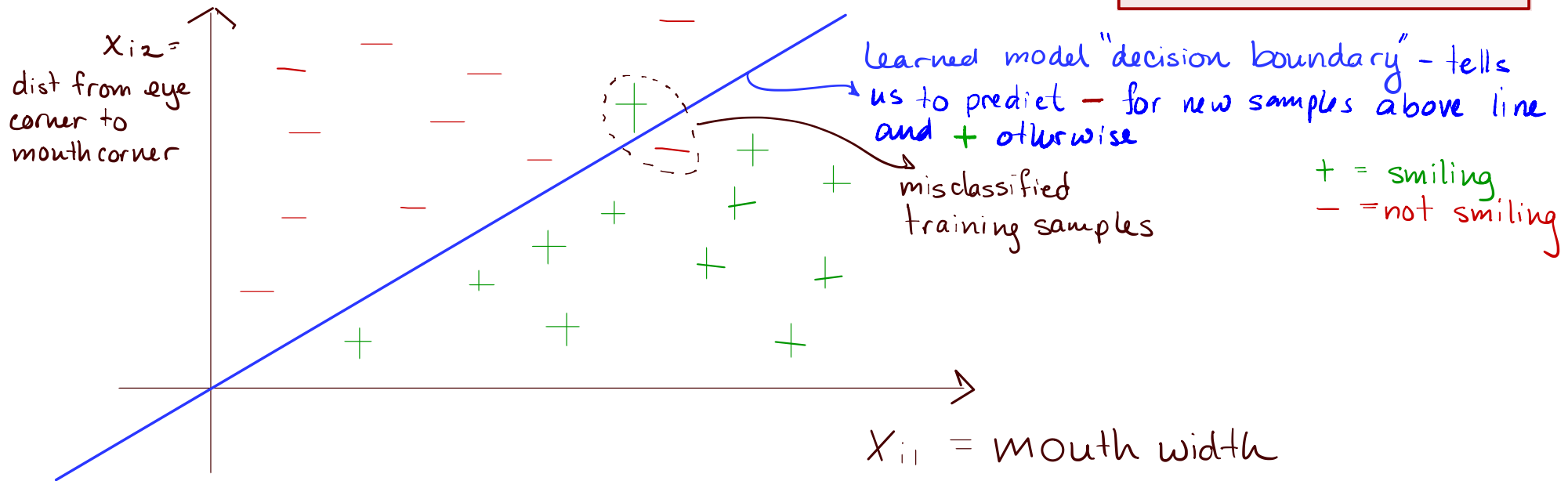
$y_i = i^{th}$ sample's label

$\underline{x}_i = i^{th}$ sample's features
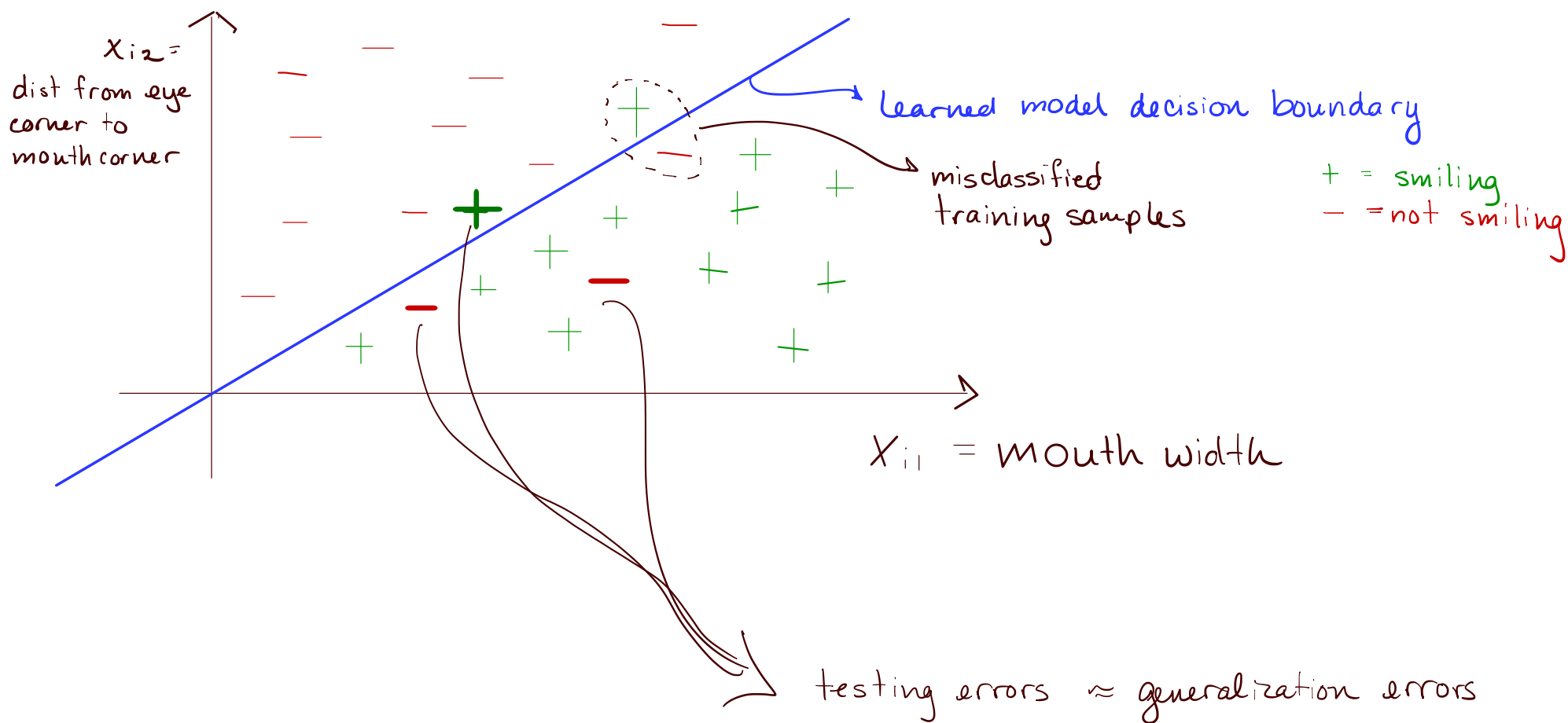
$x_{ij} = j^{th}$ feature of $i^{th}$ sample

⑤ To learn model, we choose a *loss function* = a measure of how well a model fits

data — e.g. % of samples misclassified as "smiling"

⑥ Finally, *learn the model* — search over collection of candidate models or model parameters

to find one that minimizes loss on training data

> some models learn new "features" from original input features



$X_{i2}$ =
dist from eye
corner to
mouth corner

Learned model "decision boundary" — tells us to predict — for new samples above line and + otherwise

misclassified training samples

+ = smiling
− = not smiling

$X_{i1}$ = mouth width

⑦ Characterize *generalization error* — error of our prediction on new data that was not used for training. Sometimes we estimate this using a set of test samples not used for training

$X_{i2}$ = dist from eye corner to mouth corner

learned model decision boundary

misclassified training samples

+ = smiling
− = not smiling

$X_{i1}$ = mouth width

testing errors ≈ generalization errors

# Our First ML Problem — Classification

**Learning**   we observe *training data* $(\underline{X}_i, y_i)$ for $i = 1, \ldots, n$

where $\underline{X}_i \in \mathbb{R}^p$ is a *vector* of $p$ real numbers called the *feature*

and $y_i \in \mathbb{R}$ or $y_i \in \{-1, +1\}$ or $y_i \in \{0, 1\}$ is the *label*

**Our goal**   learn a model that predicts a label $\hat{y}$ given a feature vector $\underline{X}$.

**Ex** linear model $\hat{y} = W_1 X_{o1} + W_2 X_{o2} + \quad W_p X_{op}$

$\qquad\qquad W_1, \ , W_p$ = weights to be learned from data

> This is tedious to write!
> Let's use some shorthand
> that will make many things easier

Let $\underline{w} = \begin{bmatrix} W_1 \\ W_2 \\ \vdots \\ W_p \end{bmatrix} \in \mathbb{R}^p$ be the *weight vector*
$\qquad\qquad \uparrow$ real numbers

Let $\underline{X}_o = \begin{bmatrix} X_{o1} \\ X_{o1} \\ \vdots \\ X_{op} \end{bmatrix} \in \mathbb{R}^p$ be the *feature vector*

Then our linear model can be equivalently written as

$$\hat{y} = \langle \underline{w}, \underline{x}_o \rangle = \underline{w}^T \underline{x}_o = [w_1, \quad , w_p] \begin{bmatrix} x_{o1} \\ x_{o1} \\ \vdots \\ x_{op} \end{bmatrix} = \underline{x}_o^T \underline{w} = \langle \underline{x}_o, w \rangle$$

$\underbrace{\phantom{\langle \underline{w}, \underline{x}_o \rangle}}$

Inner product of two vectors

vector transpose

---

$\boxed{\text{Ex}}$  $p = 2$  $\underline{w} = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$  When is $\hat{y} = \langle \underline{w}, \underline{x}_o \rangle > 0$ and when is $\hat{y} < 0$ ?

$\langle \underline{w}, \underline{x}_o \rangle = w_1 x_{o1} + w_2 x_{o2} > 0$

$\Rightarrow -2 x_{o1} + x_{o2} > 0$

$\Rightarrow x_{o1} < x_{o2}/2$



on this side of line, $\hat{y} > 0$

on this side of line, $\hat{y} < 0$

Line = set of $\underline{x}_o$ where $\hat{y} = 0$