# Classification: Linear Models

Ninghao Liu

University of Georgia

January 16, 2024

Some contents adopted from Stanford CS221, Percy Liang.

# Linear Classifier

**Road Map**

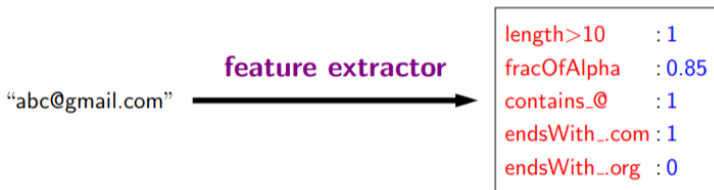Linear predictors

Learning objective
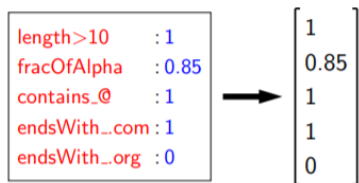
Optimization

# Feature Extraction

Example task: Predict $y$, whether a string $x$ is an email address.

Question: What properties of $x$ might be relevant for predicting $y$?

Feature extraction: Given input $x$, produce a set of (feature name, feature value) pairs.



"abc@gmail.com" → **feature extractor** →

| length>10 | : 1 |
| fracOfAlpha | : 0.85 |
| contains_@ | : 1 |
| endsWith_.com | : 1 |
| endsWith_.org | : 0 |

# Linear Classifier



$$\begin{array}{ll} \text{length>10} & : 1 \\ \text{fracOfAlpha} & : 0.85 \\ \text{contains\_@} & : 1 \\ \text{endsWith\_com} & : 1 \\ \text{endsWith\_org} & : 0 \end{array} \longrightarrow \begin{bmatrix} 1 \\ 0.85 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

For an input x, its **feature vector** is:

$$\phi(x) = [\phi(x)_1, \phi(x)_2, ..., \phi(x)_D].$$

In practice, we usually use $\boldsymbol{x}$ to denote the feature vector of input after feature extraction, i.e., $\boldsymbol{x} = \phi(x)$.

# Linear Classifier

```
length>10      :-1.2
fracOfAlpha    :0.6
contains_@     :3
endsWith_.com:2.2
endsWith_.org :1.4
...
```

For each feature $j$, have real number $w_j$ representing contribution of the feature to prediction.

# Linear Classifier



| length>10 | :-1.2 |
| fracOfAlpha | :0.6 |
| contains_@ | :3 |
| endsWith_.com | :2.2 |
| endsWith_.org | :1.4 |

| length>10 | :1 |
| fracOfAlpha | :0.85 |
| contains_@ | :1 |
| endsWith_.com | :1 |
| endsWith_.org | :0 |

Figure 1: Left: Weight vector $w$.  Right: Feature vector $\phi(x)$.

Prediction score: The weighted sum of features.

$$w^\mathsf{T} \cdot x = \sum_{j}^{D} w_j x_j \tag{1}$$

Example:

$$-1.2(1) + 0.6(0.85) + 3(1) + 2.2(1) + 1.4(0) = 4.51$$

# Linear Classifier

Weight vector $w \in \mathbb{R}^D$.

Feature vector $x \in \mathbb{R}^D$.

For binary classification, the linear classifier $f_w$ is:

$$f_w(x) = \text{sign}(w^\intercal \cdot x) = \begin{cases} +1, & \text{if } w^\intercal \cdot x > 0 \\ -1, & \text{if } w^\intercal \cdot x < 0 \\ ?, & \text{if } w^\intercal \cdot x = 0 \end{cases} \tag{2}$$

# Linear Classifier: Geometric intuition

Example:

$w = [2, -1]$

$\phi(x) = \{[2, 0], [0, 2], [2, 4]\}$

# Linear Classifier: Geometric intuition

A binary classifier $f_w$ defines a **hyperplane** with a normal vector $w$.

($\mathbb{R}^2 \rightarrow$ the hyperplane is a line; $\mathbb{R}^3 \rightarrow$ the hyperplane is a plane)

# Linear Classifier: Geometric intuition

Please note that, the complete definition of a linear classifier is:

$$f_{\mathbf{w}}(x) = sign(\mathbf{w}^{\mathsf{T}} \cdot \mathbf{x} + b), \qquad (3)$$

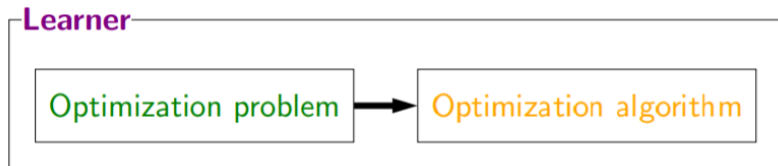where $b$ is the offset parameter.
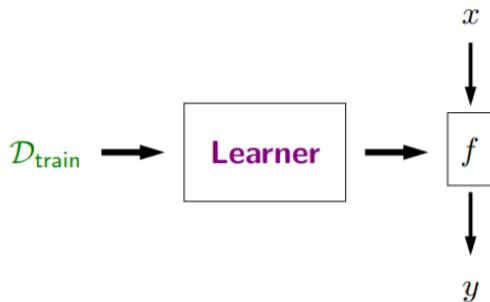
# Linear Classifier

**Road Map**

Linear predictors

Learning objective

Optimization

# Linear Classifier: Overall Framework

# Linear Classifier: Loss Function

Definition: **Loss function**.

A loss function $L(x, y, w)$ quantifies how *unhappy* you would be if you use $w$ to make a prediction on $x$ when the correct output is $y$.

Our goal is to find the best parameters $w$ that can minimize the loss function.

# Linear Classifier: Score and Margin

**Predicted** label: $y' = f_{\mathbf{w}}(x) = \text{sign}(\mathbf{w}^{\top} \cdot \mathbf{x})$
**Ground-truth** label: $y$
Example: $\mathbf{w} = [2, -1]$, $\mathbf{x} = [2, 0]$, $y = -1$

Definition: **Prediction score**.

The prediction score on an example $(x, y)$ is $\mathbf{w}^{\top} \cdot \mathbf{x}$, i.e., how confident we are in prediction.

Definition: **Margin**.

The margin on an example $(x, y)$ is $(\mathbf{w}^{\top} \cdot \mathbf{x})y$, i.e., how correct we are.

# Linear Classifier: Score and Margin

When does a binary classifier mis-classify an example?

| margin less than 0 |
|---|
| margin greater than 0 |
| score less than 0 |
| score greater than 0 |

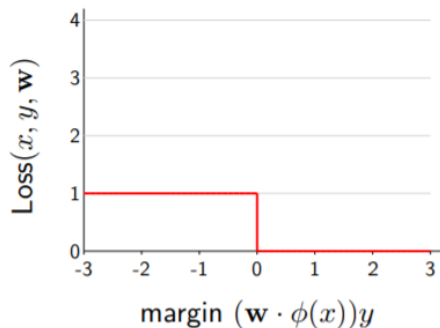# Linear Classifier: Loss function

Definition: **Zero-One loss**.

$$L_{0-1}(x, y, \boldsymbol{w}) = \mathbb{1}[y' \neq y] \tag{4}$$
$$= \mathbb{1}[f_{\boldsymbol{w}}(x) \neq y] \tag{5}$$
$$= \mathbb{1}[(\boldsymbol{w}^\mathsf{T} \cdot \boldsymbol{x})y \leq 0] \tag{6}$$
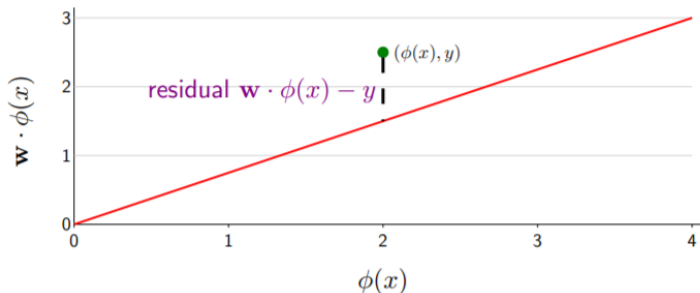
Here $\mathbb{1}$ is the indicator function, where $\mathbb{1}[True] = 1$ and $\mathbb{1}[False] = 0$.

# Linear Classifier: Loss function



$$L_{0-1}(x, y, w) = \mathbb{1}[(w^\mathsf{T} \cdot x)y \leq 0]$$

# Linear Classifier → Linear Regression



Definition: **Residual**.

The residual is $(w^\mathsf{T} \cdot x) - y$, the amount by which prediction $f_w(x) = w^\mathsf{T} \cdot x$ deviates from the target $y$.

- Define classification errors as regression errors.

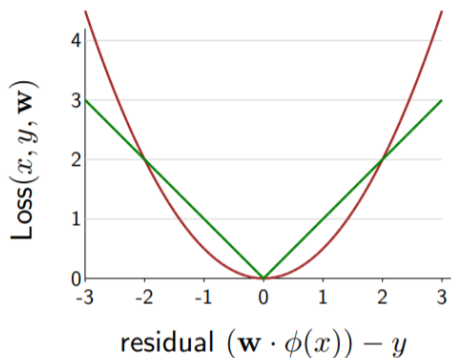# Linear Regression: Squared loss

Definition: **Squared loss**.

$$L_{squared}(x, y, \boldsymbol{w}) = (y' - y)^2 \qquad (7)$$
$$= (\boldsymbol{w}^\mathsf{T} \cdot \boldsymbol{x} - y)^2 \qquad (8)$$

Example: $\boldsymbol{w} = [2, -1]$, $\boldsymbol{x} = [2, 0]$, $y = -1$
$L_{squared}(x, y, \boldsymbol{w}) = ?$

# Linear Regression: Losses



$$L_{squared}(x, y, \boldsymbol{w}) = (\boldsymbol{w}^{\mathsf{T}} \cdot \boldsymbol{x} - y)^2 \times \frac{1}{2}$$

$$L_{absdev}(x, y, \boldsymbol{w}) = |\boldsymbol{w}^{\mathsf{T}} \cdot \boldsymbol{x} - y|$$

# Loss Minimization

So far: for one instance, $L(x, y, \boldsymbol{w})$ is easy to minimize.

How to set $\boldsymbol{w}$ to make *global* trade-offs?

- Not every instance can be happy.
- Try to make more instances happy.

Definition: **Training loss**.

$$L_{train}(\mathcal{D}_{train}, \boldsymbol{w}) = \frac{1}{|\mathcal{D}_{train}|} \sum_{(x,y) \in \mathcal{D}_{train}} L(x, y, \boldsymbol{w}) \tag{9}$$

The learning objective:

$$\min_{\boldsymbol{w}} \; L_{train}(\mathcal{D}_{train}, \boldsymbol{w}) \tag{10}$$
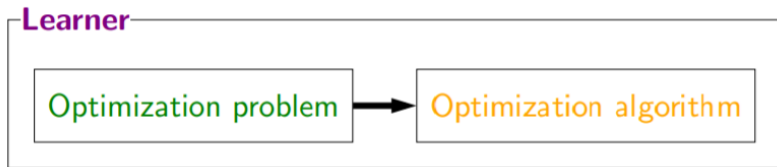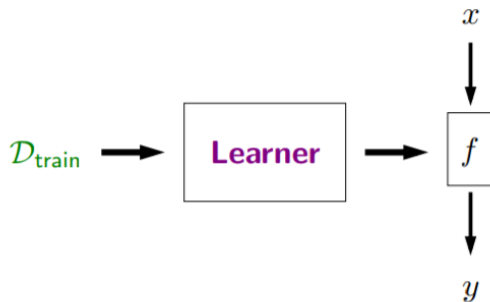
# Linear Classifier

**Road Map**

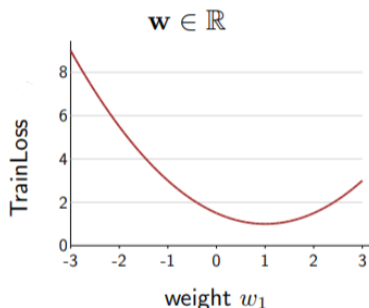Linear predictors

Learning objective

**Optimization**

# Linear Classifier: Optimization

# Linear Classifier: Optimization

The learning objective:

$$\min_{\boldsymbol{w}} \ L_{train}(\mathcal{D}_{train}, \boldsymbol{w})$$



$\mathbf{w} \in \mathbb{R}$

$\mathbf{w} \in \mathbb{R}^2$

# Linear Classifier: Optimization

The learning objective:

$$\min_{\boldsymbol{w}} \ L_{train}(\mathcal{D}_{train}, \boldsymbol{w})$$
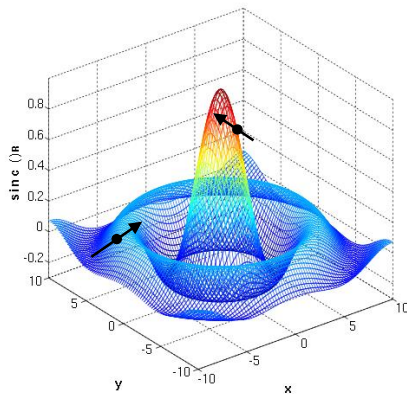
How to optimize it?

- Analytic solution (not always applicable).
- A more general solution is needed.
- **Gradient descent**.
    - An iterative algorithm.
    - At each iteration, find the fastest direction of moving $\boldsymbol{w}$ that decreases $L_{train}$.
    - Move $\boldsymbol{w}$ towards that direction.
    - Repeat the above two steps until reaching a stable $\boldsymbol{w}$.

# Gradient

Definition: **Gradient**.

The gradient $\nabla_{\boldsymbol{w}} L_{train}$ is the direction that *increases* the training loss the most.

# Gradient Descent

---

**Algorithm**  Gradient Descent

---

Initialize $\boldsymbol{w}$;
**for** $t = 1, 2, ..., T$ **do**
$\quad | \quad \boldsymbol{w} \leftarrow \boldsymbol{w} - \eta \nabla_{\boldsymbol{w}} L_{train}(\mathcal{D}_{train}, \boldsymbol{w})$
**end**

---

$\eta$: step size

$T$: number of iterations

# Gradient Descent

---

**Algorithm**  Gradient Descent

Initialize $\boldsymbol{w}$;
**for** $t = 1, 2, ..., T$ **do**
 | $\boldsymbol{w} \leftarrow \boldsymbol{w} - \eta \nabla_{\boldsymbol{w}} L_{train}(\mathcal{D}_{train}, \boldsymbol{w})$
**end**

---

How to choose the step size $\eta$?

$$0 \hspace{6cm} 1$$
$$\text{conservative, more stable} \hspace{3cm} \text{aggressive, faster} \hspace{1cm} \eta$$

Common strategies:

- Constant: e.g., $\eta = 0.01$.
- Decreasing: e.g., $\eta = 0.1/t$.

# Gradient Descent

---

**Algorithm** Gradient Descent

Initialize $\boldsymbol{w}$;
**for** $t = 1, 2, ..., T$ **do**
$\quad | \quad \boldsymbol{w} \leftarrow \boldsymbol{w} - \eta \nabla_{\boldsymbol{w}} L_{train}(\mathcal{D}_{train}, \boldsymbol{w})$
**end**

---

Example: Least squares regression

$$L_{train}(\mathcal{D}_{train}, \boldsymbol{w}) = \frac{1}{|\mathcal{D}_{train}|} \sum_{(x,y) \in \mathcal{D}_{train}} (\boldsymbol{w}^\mathsf{T} \cdot \boldsymbol{x} - y)^2$$

$$\nabla_{\boldsymbol{w}} L_{train}(\mathcal{D}_{train}, \boldsymbol{w}) = \frac{1}{|\mathcal{D}_{train}|} \sum_{(x,y) \in \mathcal{D}_{train}} 2(\boldsymbol{w}^\mathsf{T} \cdot \boldsymbol{x} - y)\boldsymbol{x}$$

# Gradient Descent

---
**Algorithm** Gradient Descent

Initialize $\boldsymbol{w}$;
**for** $t = 1, 2, ..., T$ **do**
$\quad | \quad \boldsymbol{w} \leftarrow \boldsymbol{w} - \eta \nabla_{\boldsymbol{w}} L_{train}(\mathcal{D}_{train}, \boldsymbol{w})$
**end**

---

Gradient descent is slow!

- Each iteration requires going over all training examples.
- Costly when the dataset is large (which is common nowadays).

# Stochastic Gradient Descent

---
**Algorithm**   Stochastic Gradient Descent

---
Initialize $\boldsymbol{w}$;
**for** $t = 1, 2, ..., T$ **do**
    **for** $(x, y) \in \mathcal{D}_{train}$ **do**
        $\boldsymbol{w} \leftarrow \boldsymbol{w} - \eta \nabla_{\boldsymbol{w}} L(x, y, \boldsymbol{w})$
    **end**
**end**

---

The key idea: It's not about quality, it's about quantity.

# Linear Classifier: Summary

- Linear predictor (model architecture):

$$f_{\boldsymbol{w}}(x) \text{ based on } \boldsymbol{w}^{\mathsf{T}} \cdot \phi(x).$$

- Learning objective (goal):

$$\min_{\boldsymbol{w}} L_{train}(\mathcal{D}_{train}, \boldsymbol{w}).$$

- Optimization (train the model towards the goal):

$$\boldsymbol{w} \leftarrow \boldsymbol{w} - \eta \nabla_{\boldsymbol{w}} L_{train}(\mathcal{D}_{train}, \boldsymbol{w})$$