

Linear prediction

A linear predictor for $\mathcal{X} = \mathbb{R}^d$ is a function $h : \mathbb{R}^d \rightarrow \mathbb{R}$ such that $h(\mathbf{x}) = f(\mathbf{w}^\top \mathbf{x})$ for some $\mathbf{w} \in \mathbb{R}^d$, where $f : \mathbb{R} \rightarrow \mathbb{R}$ is sometimes called the activation function. In case of regression, f is often the identity function and so $h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$. In case of binary classification, $h(\mathbf{x}) = \text{sgn}(\mathbf{w}^\top \mathbf{x} - c)$ for some $c \in \mathbb{R}$, where $\text{sgn}(z) = 1$ if $z > 0$ and -1 otherwise. We mostly focus on classification and return to regression only at the end.

Hyperplanes. Geometrically, a linear classifier defines an hyperplane $H \subset \mathbb{R}^d$ such that, if H^+, H^- are the two halfspaces on each side of the hyperplane, we have

$$h(\mathbf{x}) = \begin{cases} +1 & \text{if } \mathbf{x} \in H^+ \\ -1 & \text{if } \mathbf{x} \in H^- \end{cases}$$

Recall that an hyperplane with coefficients \mathbf{w}, c is defined by $\{\mathbf{x} \in \mathbb{R}^d : \mathbf{w}^\top \mathbf{x} = c\}$, where $\mathbf{w}^\top \mathbf{x} = \|\mathbf{w}\| \|\mathbf{x}\| \cos \theta$ and θ is the angle between \mathbf{w} and \mathbf{x} . Also, $\|\mathbf{x}\| \cos \theta$ is the length of the projection of \mathbf{x} onto \mathbf{w} and, symmetrically, $\|\mathbf{w}\| \cos \theta$ is the length of the projection of \mathbf{w} onto \mathbf{x} . Hence, the hyperplane $\{\mathbf{x} \in \mathbb{R}^d : \mathbf{w}^\top \mathbf{x} = c\}$ is orthogonal to \mathbf{w} and intersects it at distance $c/\|\mathbf{w}\|$ from the origin.

The halfspaces H^+ e H^- defined by the hyperplane $\{\mathbf{x} \in \mathbb{R}^d : \mathbf{w}^\top \mathbf{x} = c\}$ are

$$H^+ \equiv \{\mathbf{x} : \mathbf{w}^\top \mathbf{x} > c\} \quad \text{and} \quad H^- \equiv \{\mathbf{x}' : \mathbf{w}^\top \mathbf{x}' \leq c\}$$

That is, all points \mathbf{x} whose projection onto \mathbf{w} has length strictly bigger than $c/\|\mathbf{w}\|$, and all points \mathbf{x}' whose projection onto \mathbf{w} has length not larger than $c/\|\mathbf{w}\|$.

Hyperplanes of the form $\{\mathbf{x} \in \mathbb{R}^d : \mathbf{w}^\top \mathbf{x} = 0\}$ pass through the origin and are called *homogeneous*. Any non-homogeneous hyperplane $\{\mathbf{x} \in \mathbb{R}^d : \mathbf{w}^\top \mathbf{x} = c\}$, with $c \neq 0$, is equivalent to the homogeneous hyperplane $\{\mathbf{x} \in \mathbb{R}^{d+1} : \mathbf{v}^\top \mathbf{x} = 0\}$ with $\mathbf{v} = (w_1, \dots, w_d, -c)$ in the following sense: $\mathbf{w}^\top \mathbf{x} - c = \mathbf{v}^\top \mathbf{x}'$ for all $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{x}' = (x_1, \dots, x_d, 1) \in \mathbb{R}^{d+1}$. For this reason, without any loss of generality we only deal with algorithms that learn linear predictors corresponding to homogeneous hyperplanes. This amounts to saying that we automatically add an extra feature with value 1 to all of our data points.

Training linear classifiers. Recall that a linear classifier is a predictor h such that $h(\mathbf{x}) = \text{sgn}(\mathbf{w}^\top \mathbf{x})$. Clearly, $\text{sgn}(\mathbf{w}^\top \mathbf{x}) = \text{sgn}(\|\mathbf{w}\| \|\mathbf{x}\| \cos \theta) = \text{sgn}(\cos \theta)$. As the classification is only determined by the angle θ between \mathbf{w} and \mathbf{x} , the value of $\|\mathbf{w}\|$ is immaterial and we may take $\|\mathbf{w}\| = 1$. Let \mathcal{H}_d be the family of linear classifiers $h(\mathbf{x}) = \text{sgn}(\mathbf{w}^\top \mathbf{x})$ for $\mathbf{w} \in \mathbb{R}^d$ such that $\|\mathbf{w}\| = 1$. Note that the zero-one loss $\mathbb{I}\{h(\mathbf{x}_t) \neq y_t\}$ can be equivalently¹ rewritten as $\mathbb{I}\{y_t \mathbf{w}^\top \mathbf{x}_t \leq 0\}$.

¹Note that $y_t \mathbf{w}^\top \mathbf{x}_t = 0$ only when $\mathbf{w} = \mathbf{0}$ (we assume $\mathbf{x}_t \neq \mathbf{0}$ for all t). In this case, $\text{sgn}(\mathbf{w}^\top \mathbf{x}_t) = -1$ and so the classification is actually correct when $y_t = -1$. Hence, using $\mathbb{I}\{y_t \mathbf{w}^\top \mathbf{x}_t \leq 0\}$ to count mistakes we may overcount by at most one.

Consider the ERM algorithm for zero-one loss that, given a training set S containing examples $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m) \in \mathbb{R}^d \times \{-1, 1\}$, outputs

$$h_S = \operatorname{argmin}_{h \in \mathcal{H}_d} \frac{1}{m} \sum_{t=1}^m \mathbb{I}\{h(\mathbf{x}_t) \neq y_t\} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{m} \sum_{t=1}^m \mathbb{I}\{y_t \mathbf{w}^\top \mathbf{x}_t \leq 0\}. \quad (1)$$

Unfortunately, it is unlikely to find an efficient implementation of ERM for linear classifiers with zero-one loss. In fact, the decision problem associated with finding h_S is NP-complete even when $\mathbf{x}_t \in \{0, 1\}^d$ for $t = 1, \dots, m$. More precisely, introduce the following decision problem.

MinDisagreement

Instance: Pairs $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m) \in \{0, 1\}^d \times \{-1, 1\}$. Integer k .

Question: Is there $\mathbf{w} \in \mathbb{Q}^d$ such that $y_t \mathbf{w}^\top \mathbf{x}_t \leq 0$ for at most k indices $t = 1, \dots, m$?

The following result can be shown.

Theorem 1. *MinDisagreement is NP-complete.*

In addition to that, the following stronger hardness-of-approximation result can be also shown.

MinDisOpt

Instance: Pairs $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m) \in \{0, 1\}^d \times \{-1, 1\}$.

Solution: A point $\mathbf{w} \in \mathbb{Q}^d$ minimizing the number of indices $t = 1, \dots, m$ such that $y_t \mathbf{w}^\top \mathbf{x}_t \leq 0$.

Given an instance S (i.e., a training set) of MinDisOpt, let $\operatorname{Opt}(S)$ the number of examples in S that are misclassified by the ERM classifier h_S . In other words, $\operatorname{Opt}(S)/m = \ell_S(h_S)$.

Theorem 2. *If $P \neq NP$, then for all $C > 0$ there are no polynomial time algorithms that approximately solve every instance S of MinDisOpt with a number of misclassified examples bounded by $C \times \operatorname{Opt}(S)$.*

This implies that, unless $P = NP$ (which is believed unlikely), there are no efficient algorithms that approximate the solution of (1) to within any constant factor. Here efficient means with running time polynomial in the input size md .

The ERM problem (1) becomes easier when the training set is **linearly separable**. A training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$ is linearly separable where there exists a linear classifier with zero training error. In other words, there exists a separating hyperplane $\mathbf{u} \in \mathbb{R}^d$ such that

$$\gamma(\mathbf{u}) \stackrel{\text{def}}{=} \min_{t=1, \dots, m} y_t \mathbf{u}^\top \mathbf{x}_t > 0$$

The quantity $\gamma(\mathbf{u})$ is known as the **margin** of \mathbf{u} on the training set. The scaled margin $\gamma(\mathbf{u})/\|\mathbf{u}\|$ measures the distance between the separating hyperplane and the closest training example.

Now observe that the ERM problem (1) can be expressed as a system of linear inequalities,

$$y_t \mathbf{w}^\top \mathbf{x}_t > 0 \quad t = 1, \dots, m .$$

When the training set is linearly separable, the system has at least a solution. This solution can be found using an algorithm for linear programming.

We now introduce a very simple algorithm for learning linear classifiers that can be used to solve the ERM problem in the linearly separable case. The Perceptron algorithm finds a homogeneous separating hyperplane by running through the training examples one after the other. The current linear classifier is tested on each training example and, in case of misclassification, the associated hyperplane is adjusted. Note that if the algorithm terminates, then \mathbf{w} is a separating hyperplane.

```

Data: Training set  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$ 
 $\mathbf{w} = (0, \dots, 0)$ 
while true do
  | for  $t = 1, \dots, m$  do      (epoch)
  | | if  $y_t \mathbf{w}^\top \mathbf{x}_t \leq 0$  then
  | | |  $\mathbf{w} \leftarrow \mathbf{w} + y_t \mathbf{x}_t$   (update)
  | | end
  | if no update in last epoch then break
end
Output:  $\mathbf{w}$ 

```

Algorithm 1: The Perceptron algorithm (for the linearly separable case)

The update $\mathbf{w} \leftarrow \mathbf{w} + y_t \mathbf{x}_t$ when $y_t \mathbf{w}^\top \mathbf{x}_t \leq 0$ makes $y_t \mathbf{w}^\top \mathbf{x}_t$ bigger. Indeed,

$$y_t (\mathbf{w} + y_t \mathbf{x}_t)^\top \mathbf{x}_t = y_t \mathbf{w}^\top \mathbf{x}_t + \|\mathbf{x}_t\|^2 > y_t \mathbf{w}^\top \mathbf{x}_t$$

Geometrically, each update moves \mathbf{w} towards \mathbf{x}_t if $y_t = 1$ and moves \mathbf{w} away from \mathbf{x}_t if $y_t = -1$.

We now prove that Perceptron always terminates on linearly separable training sets.

Theorem 3 (Convergence of Perceptron). *Let $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$ be a linearly separable training set. Then the Perceptron algorithm terminates after a number of updates not bigger than*

$$\left(\min_{\mathbf{u}: \gamma(\mathbf{u}) \geq 1} \|\mathbf{u}\|^2 \right) \left(\max_{t=1, \dots, m} \|\mathbf{x}_t\|^2 \right) \quad (2)$$

The apparently stonger margin constraint $\gamma(\mathbf{u}) \geq 1$ is actually achievable by any separating hyperplane \mathbf{u} . Indeed, if $\gamma(\mathbf{u}) > 0$, then $y_t \mathbf{u}^\top \mathbf{x}_t \geq \gamma(\mathbf{u})$ is equivalent to $y_t \mathbf{v}^\top \mathbf{x}_t \geq 1$ for $\mathbf{v} = \mathbf{u} / \gamma(\mathbf{u})$. Hence, $\gamma(\mathbf{u}) \geq 1$ can be achieved simply by rescaling \mathbf{u} .

PROOF. Let $\mathbf{w}_0 = (0, \dots, 0)$ be the initial hyperplane. Let \mathbf{w}_M be the hyperplane after M updates and let $t_M \in \{1, \dots, m\}$ be the index of the training example $(\mathbf{x}_{t_M}, y_{t_M})$ that caused the M -th update $\mathbf{w}_M = \mathbf{w}_{M-1} + y_{t_M} \mathbf{x}_{t_M}$. We prove an upper bound on M by deriving upper and lower bounds on $\|\mathbf{w}_M\| \|\mathbf{u}\|$. We start by observing that

$$\|\mathbf{w}_M\|^2 = \|\mathbf{w}_{M-1} + y_{t_M} \mathbf{x}_{t_M}\|^2 = \|\mathbf{w}_{M-1}\|^2 + \|\mathbf{x}_{t_M}\|^2 + 2 y_{t_M} \mathbf{w}_{M-1}^\top \mathbf{x}_{t_M} \leq \|\mathbf{w}_{M-1}\|^2 + \|\mathbf{x}_{t_M}\|^2$$

because $y_{t_M} \mathbf{w}_{M-1}^\top \mathbf{x}_{t_M} \leq 0$ due to the update $\mathbf{w}_M = \mathbf{w}_{M-1} + y_{t_M} \mathbf{x}_{t_M}$. Iterating this argument M times, and recalling that $\mathbf{w}_0 = (0, \dots, 0)$, we obtain

$$\|\mathbf{w}_M\|^2 \leq \|\mathbf{w}_0\|^2 + \sum_{i=1}^M \|\mathbf{x}_{t_i}\|^2 \leq M \left(\max_{t=1, \dots, m} \|\mathbf{x}_t\|^2 \right).$$

Hence

$$\|\mathbf{w}_M\| \|\mathbf{u}\| \leq \|\mathbf{u}\| \left(\max_{t=1, \dots, m} \|\mathbf{x}_t\| \right) \sqrt{M}.$$

To prove the lower bound, fix any separating hyperplane \mathbf{u} with $\gamma(\mathbf{u}) \geq 1$ and let θ be the angle between \mathbf{u} and \mathbf{w}_M . We have

$$\begin{aligned} \|\mathbf{w}_M\| \|\mathbf{u}\| &\geq \|\mathbf{w}_M\| \|\mathbf{u}\| \cos(\theta) && \text{(since } -1 \leq \cos(\theta) \leq 1) \\ &= \mathbf{w}_M^\top \mathbf{u} && \text{(by definition of inner product)} \\ &= (\mathbf{w}_{M-1} + y_{t_M} \mathbf{x}_{t_M})^\top \mathbf{u} \\ &= \mathbf{w}_{M-1}^\top \mathbf{u} + y_{t_M} \mathbf{u}^\top \mathbf{x}_{t_M} \\ &\geq \mathbf{w}_{M-1}^\top \mathbf{u} + 1 \end{aligned}$$

where the last inequality holds because $1 \leq \gamma(\mathbf{u}) \leq y_t \mathbf{u}^\top \mathbf{x}_t$ for all $t = 1, \dots, m$. Iterating M times we get

$$\|\mathbf{w}_M\| \|\mathbf{u}\| \geq \mathbf{w}_0^\top \mathbf{u} + M = M$$

Where we used $\mathbf{w}_0^\top \mathbf{u} = 0$ since $\mathbf{w}_0 = (0, \dots, 0)$. Combining upper and lower bound we obtain

$$M \leq \|\mathbf{u}\| \left(\max_{t=1, \dots, M} \|\mathbf{x}_t\| \right) \sqrt{M}.$$

Solving for M , and recalling the choice of \mathbf{u} , we obtain (2). Hence, the update count M cannot grow larger than (2). Since the algorithm stops when no more updates are possible, we conclude that the Perceptron terminates after a bounded number of updates. \square

Note that the Perceptron convergence theorem does not imply that the Perceptron algorithm terminates in polynomial time on any linearly separable training set. Although each update takes constant time $\Theta(d)$, the number of updates can still be exponential in md whenever $\gamma(\mathbf{u}) \geq 1$ only for those \mathbf{u} whose length $\|\mathbf{u}\|$ is very big. Or, equivalently, when the margin $\gamma(\mathbf{u})$ is very small for any linear separator \mathbf{u} such that $\|\mathbf{u}\| = 1$.

Linear regression. In linear regression our predictors are linear functions $h : \mathbb{R}^d \rightarrow \mathbb{R}$ each parameterized by a vector $\mathbf{w} \in \mathbb{R}^d$ of real coefficients. That is, $h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$.

Given a training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m) \in \mathbb{R}^d \times \mathbb{R}$, the linear regression predictor is ERM with respect to the square loss,

$$\mathbf{w}_S = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \sum_{t=1}^m (\mathbf{w}^\top \mathbf{x}_t - y_t)^2$$

Now let $\mathbf{v} = (\mathbf{w}^\top \mathbf{x}_1, \dots, \mathbf{w}^\top \mathbf{x}_m)$ and $\mathbf{y} = (y_1, \dots, y_m)$. Then

$$\sum_{t=1}^m (\mathbf{w}^\top \mathbf{x}_t - y_t)^2 = \|\mathbf{v} - \mathbf{y}\|^2.$$

View all vectors as column vectors. Since $\mathbf{v} = S\mathbf{w}$, where S is the $m \times d$ **design matrix** such that $S^\top = [\mathbf{x}_1, \dots, \mathbf{x}_m]$, we may also write

$$\mathbf{w}_S = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \|S\mathbf{w} - \mathbf{y}\|^2 .$$

Since $F(\mathbf{w}) = \|S\mathbf{w} - \mathbf{y}\|^2$ is a convex function, the minimizer satisfies the condition $\nabla F(\mathbf{w}) = \mathbf{0}$.

Using matrix calculus, we have that $\nabla \|S\mathbf{w} - \mathbf{y}\|^2 = 2S^\top(S\mathbf{w} - \mathbf{y})$. Hence, $\nabla \|S\mathbf{w} - \mathbf{y}\|^2 = \mathbf{0}$ for $\mathbf{w} = (S^\top S)^{-1}S^\top \mathbf{y}$ provided $S^\top S$ is nonsingular (i.e., invertible) —which is equivalent to $\mathbf{x}_1, \dots, \mathbf{x}_m$ spanning \mathbb{R}^d . When this happens, we have that the ERM with respect to the square loss is $\hat{\mathbf{w}} = (S^\top S)^{-1}S^\top \mathbf{y}$.

Ridge Regression. When $S^\top S$ is nearly singular, $\hat{\mathbf{w}}$ is highly sensitive to perturbations of the training set. This instability increases the estimation error (or variance). A more stable predictor is obtained by introducing a **regularizer** in the ERM functional which increases the approximation error (or bias) and reduces the variance with a beneficial effect on the risk.

In other words, instead of defining \mathbf{w}_S by

$$\mathbf{w}_S = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \|S\mathbf{w} - \mathbf{y}\|^2$$

we use the regularized form, also known as Ridge Regression,

$$\mathbf{w}_{S,\alpha} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \|S\mathbf{w} - \mathbf{y}\|^2 + \alpha \|\mathbf{w}\|^2$$

where $\alpha > 0$ is the regularization parameter. When $\alpha \rightarrow 0$ we recover the standard linear regression solution. When $\alpha \rightarrow \infty$, the solution $\hat{\mathbf{w}}_{S,\alpha}$ becomes the zero vector. This shows that α can be used to control the bias of the algorithm.

Similarly to before, we have that

$$\nabla \left(\|S\mathbf{w} - \mathbf{y}\|^2 + \alpha \|\mathbf{w}\|^2 \right) = 2S^\top(S\mathbf{w} - \mathbf{y}) + 2\alpha\mathbf{w} .$$

Hence, the gradient vanishes for $\mathbf{w} = \hat{\mathbf{w}}_{S,\alpha} = (\alpha I + S^\top S)^{-1}S^\top \mathbf{y}$. Note that we do not have to worry anymore about the singularity of $S^\top S$. Indeed, if $\lambda_1 \geq \dots \geq \lambda_d \geq 0$ are the eigenvalues of $S^\top S$, the eigenvalues of $\alpha I + S^\top S$ are simply $\alpha + \lambda_1 \geq \dots \geq \alpha + \lambda_d > 0$. Hence, $\alpha I + S^\top S$ is invertible for all $\alpha > 0$.