

SGD Algorithm:

$$\vec{w}^{(0)} \leftarrow \vec{0};$$

for $t \leftarrow 0$ to $T-1$ do {

pick i uniformly at random from $\{1, \dots, m\}$;

$$\vec{w}^{(t+1)} \leftarrow \vec{w}^{(t)} - \eta \nabla \ell(\vec{w}^{(t)}, (\vec{x}_i, \tilde{y}_i));$$

}

$$\text{return } \bar{w} = \frac{1}{T} \sum_{t=1}^T \vec{w}^{(t)};$$

To simplify the notation, we are going to compute $\nabla \ell(\vec{w}^{(t)}, (\vec{x}_i, \tilde{y}_i))$:

$$\nabla \ell(\vec{w}^{(t)}, (\vec{x}_i, \tilde{y}_i)) = \begin{cases} 0 & \text{if } y_i \langle \vec{w}, \vec{x}_i \rangle > 0 \rightarrow y_i \text{ correctly classifies } x_i \\ \nabla (-y_i \langle \vec{w}, \vec{x}_i \rangle) & \text{otherwise} \end{cases}$$

Assume that $y_i \langle \vec{w}, \vec{x}_i \rangle < 0$:

$$\nabla (-y_i \langle \vec{w}, \vec{x}_i \rangle) = \left[\frac{\partial (-y_i \langle \vec{w}, \vec{x}_i \rangle)}{\partial w_1}, \dots, \frac{\partial (-y_i \langle \vec{w}, \vec{x}_i \rangle)}{\partial w_d} \right]^T$$

$$\text{let } \vec{x}_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{id} \end{bmatrix}. \text{ Since } -y_i \langle \vec{w}, \vec{x}_i \rangle = -y_i \sum_{j=1}^d (w_j x_{ij})$$

$$\frac{\partial (-y_i \langle \vec{w}, \vec{x}_i \rangle)}{\partial w_j} = -y_i x_{ij}$$

$$\nabla \ell(w, (\vec{x}_i, \vec{y}_i)) = \begin{bmatrix} -y_i x_{i1} & \dots & -y_i x_{id} \end{bmatrix}^T$$

$$= -y_i \vec{x}_i$$

Then we can rewrite the pseudo code:

SGD Algorithm:

$$\vec{w}^{(0)} \leftarrow \vec{0};$$

for $t = 0$ to $T-1$ do {

 pick i uniformly at random from $\{1, \dots, m\}$;

 if $y_i \langle \vec{w}^{(t)}, \vec{x}_i \rangle < 0$ then {

$$\vec{w}^{(t+1)} \leftarrow \vec{w}^{(t)} + \eta y_i \vec{x}_i;$$

 }

}

$$\text{return } \bar{w} = \frac{1}{T} \sum_{t=1}^T \vec{w}^{(t)};$$

Comparison:

perceptron	SGD perceptron
1) Choose a misclassified point	choose a point at random → the main difference
2) $\eta = 1$	η is a parameter → just a generalization
3) return "best" \vec{w}^*	return \vec{w} → implementation choice (average/best/...)

We can speed up the SGD perceptron, at each iteration, by picking a misclassified point at random



SGD perceptron is the perceptron

Linear Regression

$$\mathcal{X} = \mathbb{R}^d, \mathcal{Y} = \mathbb{R}$$

↳ Regression with linear models (wow!)

Hypothesis class:

$$\mathcal{H}_{reg} = L_d = \{\mathbf{x} \mapsto \langle \mathbf{w}, \mathbf{x} \rangle + b : \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}\}$$

↳ Linear

↳ This is already a real value

Note: $h \in \mathcal{H}_{reg} : \mathbb{R}^d \rightarrow \mathbb{R}$

Commonly used loss function: squared-loss

$$\ell(h, (\mathbf{x}, y)) \stackrel{\text{def}}{=} (h(\mathbf{x}) - y)^2$$

Annotations: $\mathbf{x} \in \mathcal{X}$, $y \in \mathcal{Y}$, $h \in \mathcal{H}_{reg}$

⇒ empirical risk function (training error): Mean Squared Error

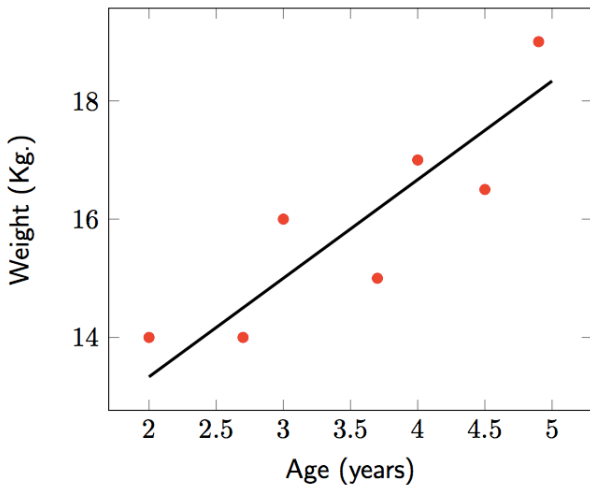
$$S = \{(\vec{x}_1, y_1), \dots, (\vec{x}_m, y_m)\} \quad L_S(h) = \frac{1}{m} \sum_{i=1}^m (h(\mathbf{x}_i) - y_i)^2$$

Annotation: $S \in \mathcal{S}$

Linear Regression - Example

$d = 1$

$$h_{\vec{w},b}(\vec{x}) = b + w \cdot x_1 \rightarrow \text{find the best line}$$



Least Squares

How to find a ERM hypothesis? *Least Squares* algorithm

Best hypothesis:

$$\arg \min_{\mathbf{w}} L_S(h_{\mathbf{w}}) = \arg \min_{\mathbf{w}} \frac{1}{m} \sum_{i=1}^m (\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i)^2$$

Equivalent formulation: \mathbf{w} minimizing *Residual Sum of Squares* (RSS), i.e.

$$\arg \min_{\mathbf{w}} \sum_{i=1}^m (\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i)^2$$

$\mathbf{w} = [b, w_1, \dots, w_d]^T$

RSS: Matrix Form

Let

$$\mathbf{X} = \begin{bmatrix} \cdots & \mathbf{x}_1 & \cdots \\ \cdots & \mathbf{x}_2 & \cdots \\ \cdots & \vdots & \cdots \\ \cdots & \mathbf{x}_m & \cdots \end{bmatrix} \rightarrow \text{each row is an instance in the training set}$$

\mathbf{X} : design matrix

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \rightarrow \text{observations from the training set}$$

\Rightarrow we have that RSS is

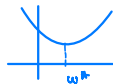
$$\sum_{i=1}^m (\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i)^2 = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$

Want to find \mathbf{w} that minimizes RSS (=objective function):

$$\arg \min_{\mathbf{w}} RSS(\mathbf{w}) = \arg \min_{\mathbf{w}} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$

↳ parabola ($x^2w^2 - 2yXw + y^2$)

How?



Compute gradient $\frac{\partial RSS(\mathbf{w})}{\partial \mathbf{w}}$ of objective function w.r.t \mathbf{w} and compare it to 0.

$$\frac{\partial RSS(\mathbf{w})}{\partial \mathbf{w}} = -2\mathbf{X}^T(\mathbf{y} - \mathbf{X}\mathbf{w})$$

Then we need to find \mathbf{w} such that

$$-2\mathbf{X}^T(\mathbf{y} - \mathbf{X}\mathbf{w}) = 0$$

$$-2\mathbf{X}^T(\mathbf{y} - \mathbf{X}\mathbf{w}) = 0$$

is equivalent to

$$\mathbf{X}^T\mathbf{X}\mathbf{w} = \mathbf{X}^T\mathbf{y}$$

If $\mathbf{X}^T\mathbf{X}$ is invertible \Rightarrow solution to ERM problem is:

$$\mathbf{w} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

Complexity Considerations

We need to compute

$$(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Algorithm:

- 1 compute $\mathbf{X}^T \mathbf{X}$: product of $(d+1) \times m$ matrix and $m \times (d+1)$ matrix
- 2 compute $(\mathbf{X}^T \mathbf{X})^{-1}$ inversion of $(d+1) \times (d+1)$ matrix
- 3 compute $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$: product of $(d+1) \times (d+1)$ matrix and $(d+1) \times m$ matrix
- 4 compute $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$: product of $(d+1) \times m$ matrix and $m \times 1$ matrix

Most expensive operation? Inversion! *→ Good since d is small compared to m*

⇒ done for $(d+1) \times (d+1)$ matrix

$\mathbf{X}^T \mathbf{X}$ not invertible?

How do we get \mathbf{w} such that

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

if $\mathbf{X}^T \mathbf{X}$ is not invertible?

Let

$$\mathbf{A} = \mathbf{X}^T \mathbf{X}$$

Let \mathbf{A}^+ be the *generalized inverse* of \mathbf{A} , i.e.:

$$\mathbf{A} \mathbf{A}^+ \mathbf{A} = \mathbf{A}$$

Proposition

If $\mathbf{A} = \mathbf{X}^T \mathbf{X}$ is not invertible, then $\hat{\mathbf{w}} = \mathbf{A}^+ \mathbf{X}^T \mathbf{y}$ is a solution to $\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$.

Computing the Generalized Inverse of \mathbf{A}

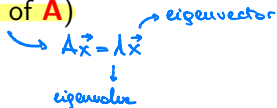
Note $\mathbf{A} = \mathbf{X}^T \mathbf{X}$ is symmetric \Rightarrow eigenvalue decomposition of \mathbf{A} :

$$\mathbf{A} = \mathbf{V} \mathbf{D} \mathbf{V}^T$$

with

- \mathbf{D} : diagonal matrix (entries = eigenvalues of \mathbf{A})
- \mathbf{V} : orthonormal matrix ($\mathbf{V}^T \mathbf{V} = \mathbf{I}_{d \times d}$)


$$\mathbf{D} = \begin{bmatrix} D_{11} & & 0 \\ & \ddots & \\ 0 & & \ddots \end{bmatrix}$$


$$\mathbf{A} \vec{x} = \lambda \vec{x}$$

eigenvector

eigenvalue

Define \mathbf{D}^+ diagonal matrix such that:

$$\mathbf{D}_{i,i}^+ = \begin{cases} 0 & \text{if } \mathbf{D}_{i,i} = 0 \\ \frac{1}{\mathbf{D}_{i,i}} & \text{otherwise} \end{cases} = \begin{bmatrix} \frac{1}{D_{11}} & & 0 \\ & \ddots & \\ 0 & & \ddots \end{bmatrix}$$

$$\mathbf{D} \mathbf{D}^+ = \begin{bmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & \ddots \end{bmatrix}$$

Let $\mathbf{A}^+ = \mathbf{V}\mathbf{D}^+\mathbf{V}^T$

Then

$$\begin{aligned}\mathbf{A}\mathbf{A}^+\mathbf{A} &= \mathbf{V}\overbrace{\mathbf{D}\mathbf{V}^T\mathbf{V}}^{\mathbf{I}}\mathbf{D}^+\mathbf{V}^T\overbrace{\mathbf{V}\mathbf{V}^T\mathbf{V}}^{\mathbf{I}}\mathbf{V}^T \\ &= \mathbf{V}\overbrace{\mathbf{D}\mathbf{D}^+\mathbf{D}}^{\mathbf{D}}\mathbf{V}^T \\ &= \mathbf{V}\mathbf{D}\mathbf{V}^T \\ &= \mathbf{A}\end{aligned}$$

$\Rightarrow \mathbf{A}^+$ is a generalized inverse of \mathbf{A} .

In practice: the Moore-Penrose generalized inverse \mathbf{A}^\dagger of \mathbf{A} is used, since it can be efficiently computed from the Singular Value Decomposition of \mathbf{A} .

Exercise

Consider a linear regression problem, where $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{Y} = \mathbb{R}$, with mean squared loss. The hypothesis set is the set of constant functions, that is $\mathcal{H} = \{h_a : a \in \mathbb{R}\}$, where $h_a(\mathbf{x}) = a$. Let $S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m))$ denote the training set.

- Derive the hypothesis $h \in \mathcal{H}$ that minimizes the training error.
- Use the result above to explain why, for a given hypothesis \hat{h} from the set of all linear models, the coefficient of determination $R^2 = 1 - \frac{\sum_{i=1}^m (\hat{h}(\mathbf{x}_i) - y_i)^2}{\sum_{i=1}^m (y_i - \bar{y})^2}$ where \bar{y} is the average of the $y_i, i = 1, \dots, m$ is a measure of how well \hat{h} performs (on the training set).

Cerco gli altri es su steam.