

1. Linearization to help classification

As was discussed in lecture, the naive straightforward way of picking the decision boundary (by looking at the mean example of each category and drawing the perpendicular bisector) is not always the best.

Consider trying to classify a set of measurements \vec{x}_i with given labels ℓ_i . For the binary case of interest here, we will think of the labels as being “+” and “-” and fold our threshold implicitly into the weights by augmenting the constant “1” in the first position of each \vec{x}_i . Now, the classification rule becomes simple. **We want to learn a vector of weights \vec{w} so that we can deem any point with $\vec{x}_i^T \vec{w} > 0$ as being a member of the “+” category and anything with $\vec{x}_i^T \vec{w} < 0$ as being a member of the “-” category.**

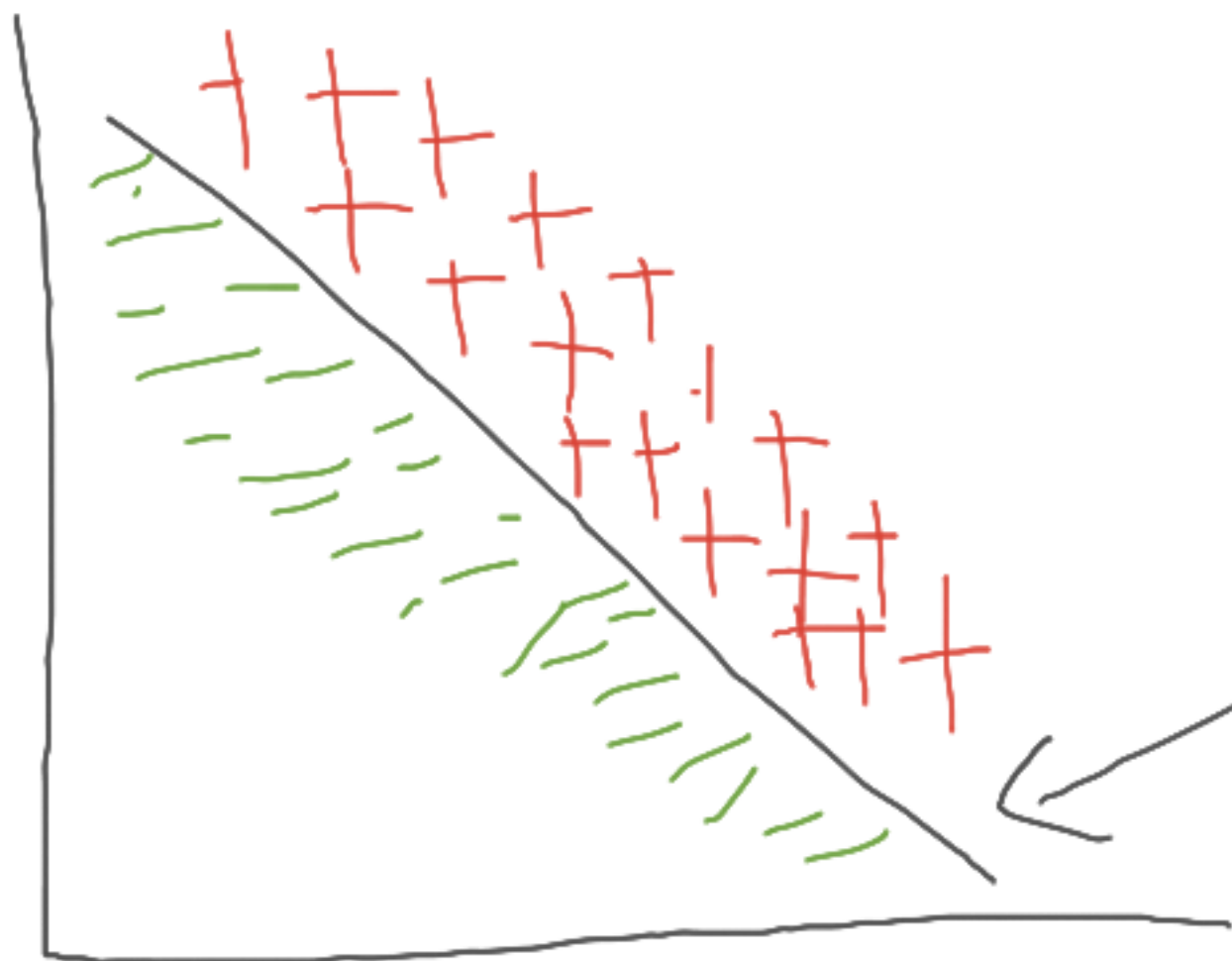
The way that we will do this, is to do a minimization in the spirit of least squares. Except, instead of necessarily using some sort of squared loss function, we will just consider a generic cost function that can depend on the label and the prediction for the point. For the i -th data point in our training data, we will incur a cost $c(\vec{x}_i^T \vec{w}, \ell_i)$ for a total cost that we want to minimize as:

$$\arg \min_{\vec{w}} c_{total}(\vec{w}) = \sum_{i=1}^m c(\vec{x}_i^T \vec{w}, \ell_i) \quad \leftarrow \text{(1) How?}$$

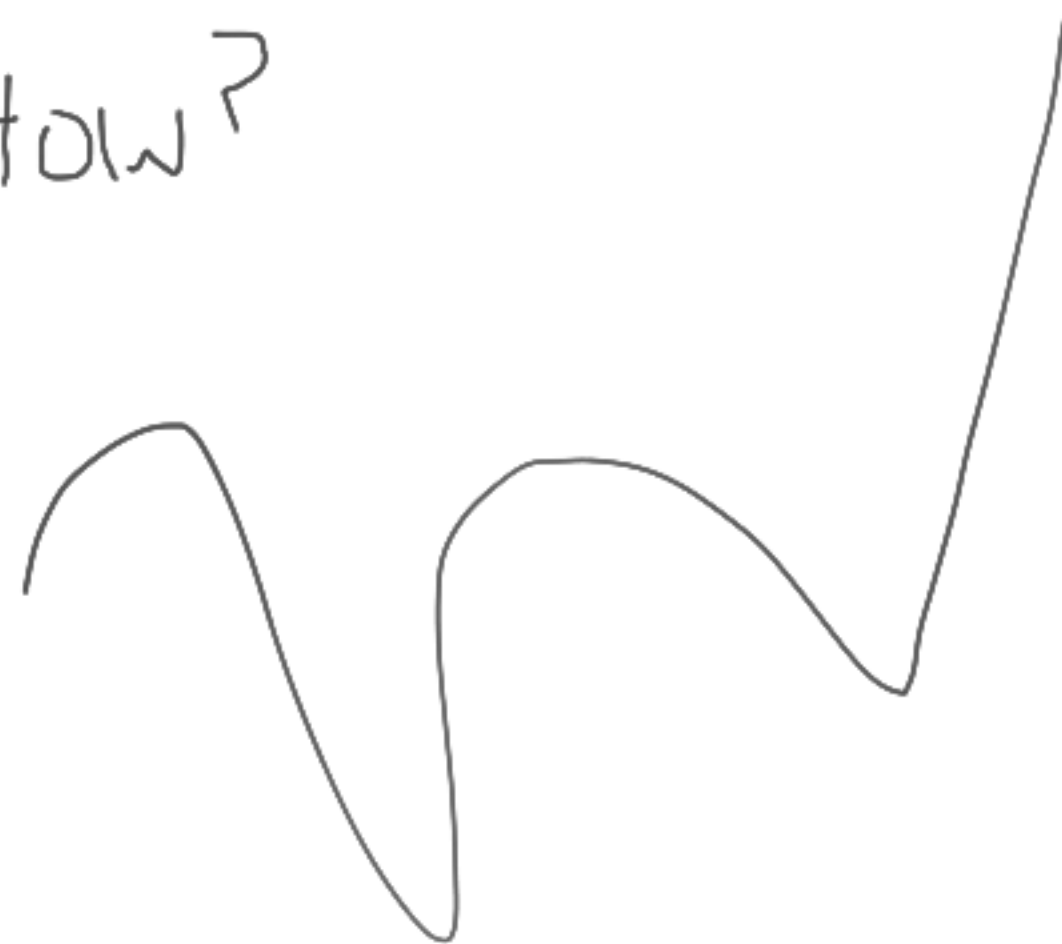
Goal: Learn \vec{w} s.t.

$$X_i \in X^+ \Rightarrow \vec{x}_i^T \vec{w} > 0$$

$$X_i \in X^- \Rightarrow \vec{x}_i^T \vec{w} < 0$$



\vec{w}
decision
boundary



Because this can be a nonlinear function, our goal is to solve this iteratively as a sequence of least-squares problems that we know how to solve.

Consider the following algorithm:

- 1: $\vec{w} = \vec{0}$ ▷ Initialize the weights to $\vec{0}$
- 2: **while** Not done **do** ▷ Iterate towards solution
- 3: Compute $\vec{w}^\top \vec{x}_i$ ▷ Generate current estimated labels
- 4: Compute $\frac{d}{d\vec{w}} c(\vec{w}^\top \vec{x}_i, \ell_i)$ ▷ Generate derivatives with respect to \vec{w} of the cost for update step
- 5: Compute $\frac{d^2}{d\vec{w}^2} c(\vec{w}^\top \vec{x}_i, \ell_i)$ ▷ Generate second derivatives of the cost for update step
- 6: $\delta\vec{w} = \text{LeastSquares}(\cdot, \cdot)$ ▷ We will derive what to call least squares on
- 7: $\vec{w} = \vec{w} + \delta\vec{w}$ ▷ Update parameters
- 8: **end while**
- 9: **Return** \vec{w}

The key step above is figuring out with what arguments to call LeastSquares while only having the labels ℓ_i and the points \vec{x}_i .

(a) Now, suppose we wanted to approximate the cost for each data point

$$c_i(\vec{w}) = c(\vec{x}_i^T \vec{w}, \ell_i)$$

where

$$\vec{w} = \begin{bmatrix} w[1] \\ \vdots \\ w[n] \end{bmatrix}$$

in the neighborhood of a weight vector \vec{w}_* . Our goal is to write out the first-order expression for approximating the cost function $c_i(\vec{w}_* + \delta\vec{w})$. This should be something in vector/matrix form like you have seen for the approximation of nonlinear systems by linear systems. We don't want to take any second derivatives just yet — only first derivatives. We have outlined a skeleton for the derivation with some parts missing. Follow the guidelines in each sub-section.

i) Comparing to eq. (2), we know that $c_i(\vec{w}_* + \delta\vec{w}) \approx c_i(\vec{w}_*) + D_{\vec{w}}c_i(\vec{w}_*)\delta\vec{w}$. **Write out the vector form of $D_{\vec{w}}c_i(\vec{w}_*)$ in terms of the partial derivatives.**

ii) **Write out the partial derivatives of $c_i(\vec{w})$ with respect to $w[g]$, the g^{th} component of \vec{w} .** (HINT: Use the linearity of derivatives and sums to compute the partial derivatives with respect to each of the $w[g]$ terms. Don't forget the chain rule and the fact that $\vec{x}_i^T \vec{w} = \sum_{j=1}^n x_i[j]w[j] = x_i[g]w[g] + \sum_{j \neq g} x_i[j]w[j]$.)

iii) With what you had above, **can you fill in the missing part to express the row vector $D_{\vec{w}}c_i(\vec{w})$?**

(3)

$$i) \frac{\partial c_i}{\partial \vec{w}} = \left[\frac{\partial c_i}{\partial w_1} \quad \dots \quad \frac{\partial c_i}{\partial w_n} \right]$$

$$iii) c'(\vec{x}_i^T \vec{w}, \ell_i) \begin{bmatrix} x_i[1] & \dots & x_i[n] \end{bmatrix} \\ = c'(\vec{x}_i^T \vec{w}, \ell_i) \vec{x}_i^T$$

$$ii) \frac{\partial c_i(\vec{w})}{\partial w_g} = \frac{\partial c(\vec{x}_i^T \vec{w}, \ell_i)}{\partial w_g} = \frac{\partial c(\vec{x}_i^T \vec{w}, \ell_i)}{\partial \vec{x}_i^T \vec{w}} \frac{\partial \vec{x}_i^T \vec{w}}{\partial w_g} \\ = c'(\vec{x}_i^T \vec{w}, \ell_i) x_i[g]$$

(b) Now, we want a better approximation that includes second derivatives. For a general function, we would look for

$$f(\vec{x}_0 + \delta\vec{x}) \approx f(\vec{x}_0) + f'(\vec{x}_0)\delta\vec{x} + \frac{1}{2}\delta\vec{x}^\top f''(\vec{x}_0)\delta\vec{x} \quad (4)$$

where $f'(\vec{x}_0)$ is an appropriate row vector and, as you've seen in the note, $f''(\vec{x}_0)$ is called Hessian that represents the second derivatives.

i) Comparing to eq. (4), we know that

$$c_i(\vec{w}_* + \delta\vec{w}) \approx c_i(\vec{w}_*) + D_{\vec{w}}c_i(\vec{w}_*)\delta\vec{w} + \frac{1}{2}\delta\vec{w}^\top H_{\vec{w}}c_i(\vec{w}_*)\delta\vec{w}$$

Write out the matrix form of $H_{\vec{w}}c_i(\vec{w}_*)$.

$$\frac{\partial^2 c_i}{\partial \vec{w}^2} = \begin{bmatrix} \frac{\partial^2 c_i}{\partial w_1^2} & \dots & \frac{\partial^2 c_i}{\partial w_1 \partial w_n} \\ \vdots & & \vdots \\ \frac{\partial^2 c_i}{\partial w_n \partial w_1} & \dots & \frac{\partial^2 c_i}{\partial w_n \partial w_n} \end{bmatrix}$$

$n \times n - |c_n(\vec{w})|$

ii) Take the second derivatives of the cost $c_i(\vec{w})$, i.e. solve for $\frac{\partial^2 c_i(\vec{w})}{\partial w[g] \partial w[h]}$.

(HINT: You should use the answer to part (a) and just take another derivative. Once again, use the linearity of derivatives and sums to compute the partial derivatives with respect to each of the $w[h]$ terms. This will give you $\frac{\partial^2}{\partial w[g] \partial w[h]}$. Don't forget the chain rule and again use the fact that $\vec{x}_i^T \vec{w} = \sum_{j=1}^n x_i[j]w[j] = x_i[h]w[h] + \sum_{j \neq h} x_i[j]w[j]$.)

$$\frac{\partial^2 c_i(\vec{w})}{\partial w[h] \partial w[g]} = \frac{\partial}{\partial w[h]} \left(\frac{\partial}{\partial w[g]} c_i(\vec{w}) \right) = \frac{\partial}{\partial w[h]} [c'(\vec{x}_i^T \vec{w}, l_i) \cdot x_i[g]]$$

$$\frac{\partial [c'(\vec{x}_i^T \vec{w}, l_i) \cdot x_i[g]]}{\partial \vec{x}_i^T \vec{w}} \cdot \frac{\partial \vec{x}_i^T \vec{w}}{\partial w[h]} = c''(\vec{x}_i^T \vec{w}, l_i) \cdot x_i[g] \cdot x_i[h]$$

$$c'(f(\vec{w}), l_i), \quad f(\vec{w}) = \vec{x}_i^T \vec{w}$$

iii) The expression in part (ii) is for the $[g, h]$ -th component of the second derivative. $\frac{1}{2}$ times this times $\delta\vec{w}[g]$ times $\delta\vec{w}[h]$ would give us that component's contribution to the second-derivative term in the approximation, and we have to sum this up over all g and h to get the total contribution of the second-derivative term in the approximation. Now, we want to group terms to restructure this into matrix-vector form by utilizing the outer-product form of matrix multiplication. **What should the space in the following expression be filled with?**

$$H_{\vec{w}}c_i(\vec{w}) = c''(\vec{x}_i^T \vec{w}, l_i) \underline{\quad} \vec{x}_i \vec{x}_i^T$$

$$c''(\vec{x}_i^T \vec{w}, l_i) \begin{bmatrix} x_i[1]x_i[1] & \dots & x_i[1]x_i[n] \\ \vdots & & \vdots \\ x_i[1]x_i[n] & \dots & x_i[n]x_i[n] \end{bmatrix} = c''(\vec{x}_i^T \vec{w}, l_i) \vec{x}_i \vec{x}_i^T$$

(c) Now we have successfully expressed the second order approximation of $c_i(\vec{w}_* + \delta \vec{w})$. Since we eventually want to minimize the total cost $c_{total}(\vec{w}) = \sum_{i=1}^m c_i(\vec{w})$, can you write out the second order approximation of $c_{total}(\vec{w}_* + \delta \vec{w})$ using results from (a) and (b)?

$$\begin{aligned}
 c_i(\vec{w}_* + \delta \vec{w}) &\approx c_i(\vec{w}_*) + \frac{\partial c_i(\vec{w})}{\partial \vec{w}}(\vec{w}_*) \cdot \delta \vec{w} + \frac{1}{2} \delta \vec{w}^T \frac{\partial^2 c_i(\vec{w})}{\partial \vec{w}^2}(\vec{w}_*) \delta \vec{w} \\
 &= c_i(\vec{w}_*) + c'(\vec{x}_i^T \vec{w}_*, l_i) \vec{x}_i^T \delta \vec{w} + \frac{1}{2} \delta \vec{w}^T c''(\vec{x}_i^T \vec{w}_*, l_i) \vec{x}_i \vec{x}_i^T \delta \vec{w}
 \end{aligned}$$

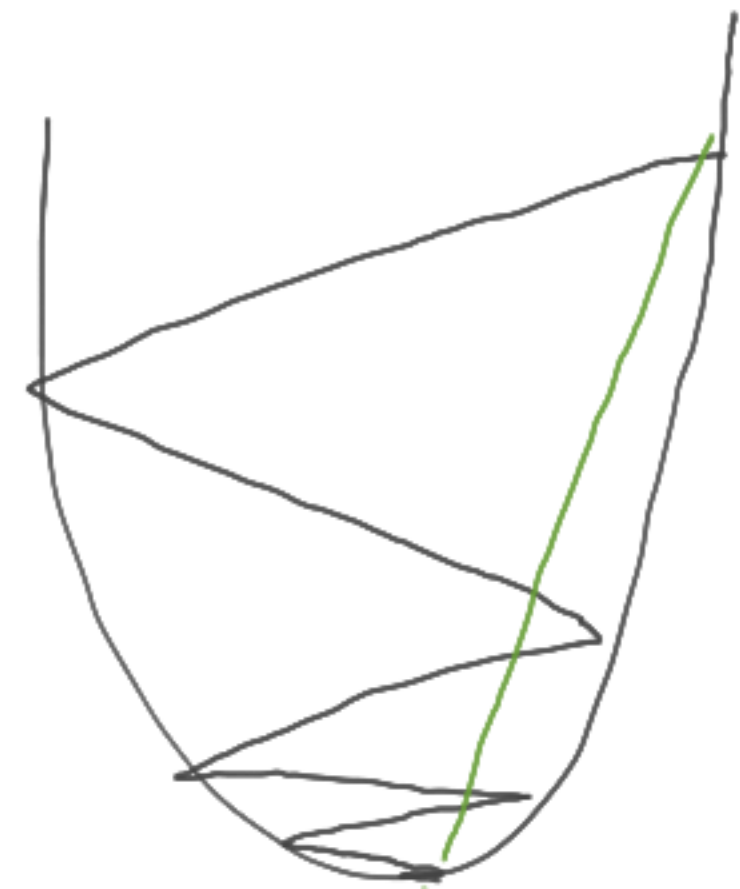
$$c_{total}(\vec{w}_* + \delta \vec{w}) = \sum_{i=1}^m c_i(\vec{w}_* + \delta \vec{w})$$

(d) In this part we explore "Newton's method", as an alternative to the iterative least squares formulation above. Recall that for a differentiable function $f(\vec{w})$, we can set $D_{\vec{w}}f(\vec{w}^*) = \vec{0}^\top$ to find a minimum. Consider a linearization of the derivative of our cost function, $D_{\vec{w}}c^\top$ at the point \vec{w} . **Using the linearization for $D_{\vec{w}}c(\vec{w})^\top$ and the fact that $D_{\vec{w}}c|_{\vec{w}^*} = \vec{0}^\top$ to derive an update.**

$$\frac{\partial c_i(\vec{w}_1)}{\partial \vec{w}} = \frac{\partial c_i(\vec{w}_0)}{\partial \vec{w}} + \frac{\partial^2 c_i(\vec{w}_0)}{\partial \vec{w}^2} (\vec{w}_1 - \vec{w}_0)$$

$$0 = \frac{\partial^2 c_i(\vec{w}_0)}{\partial \vec{w}^2} (\vec{w}_1 - \vec{w}_0) + \frac{\partial c_i(\vec{w}_0)}{\partial \vec{w}}$$

$$\vec{w}_1 = \vec{w}_0 - \left(\frac{\partial^2 c_i(\vec{w}_0)}{\partial \vec{w}^2} \right)^{-1} \frac{\partial c_i(\vec{w}_0)}{\partial \vec{w}}$$





$$\vec{x}_i^T \vec{w} < 0$$



$$\vec{x}_i^T \vec{w} > 0$$

Feedback: <https://tinyurl.com/manav16b>